

Hierarchical Stochastic Model in Bayesian Inference: Theoretical Implications and Efficient Approximation*

S. Wu^{†1}, P. Angelikopoulos¹, J. L. Beck² and P. Koumoutsakos¹

¹Computational Science and Engineering Laboratory,
Claussiusstrasse 33, ETH-Zurich 8092, Switzerland.

²Department of Mechanical and Civil Engineering, California
Institute of Technology, Pasadena, CA 91102, USA.

November 10, 2016

Abstract

We classify two types of Hierarchical Bayesian Model found in the literature as Hierarchical Prior Model (HPM) and Hierarchical Stochastic Model (HSM). Then, we focus on studying the theoretical implications of the HSM. Using examples of polynomial functions, we show that the HSM is capable of separating different types of uncertainties in a system and quantifying uncertainty of reduced order models under the Bayesian model class selection framework. To tackle the huge computational cost for analyzing HSM, we propose an efficient approximation scheme based on Importance Sampling and Empirical Interpolation Method. We illustrate our method using two examples — a Molecular Dynamics simulation for Krypton and a pharmacokinetic/pharmacodynamic model for cancer drug.

Key words: Hierarchical Bayesian, Importance Sampling, Empirical Interpolation Method, Molecular Dynamics, Pharmacokinetics

1 Introduction

Hierarchical Bayesian Model (HBM) is a powerful modeling tool extended from the classical Bayesian framework. It has been used to solve many difficult practical problems, such as modeling heterogeneous data, calibrating system with

*This work was supported by the European Research Council (Advanced Investigator Award no. 341117).

[†]Currently at The Institute of Statistical Mathematics, Japan.

multiple objectives and inducing sparsity in a model. Its applications span across various fields, such as, social science, economics, physics, medical science, civil engineering, etc [28, 9, 2, 14].

In a classical Bayesian setting, users define a stochastic model with parameters $\vec{\theta}$ for a forward problem that predicts a quantity of interest and a prior distribution of $\vec{\theta}$. When data is available, the Bayes' Theorem is used to solve the inverse problem by finding the posterior distribution of $\vec{\theta}$. For a complex system in reality, we further parameterize the stochastic model and the prior with hyperparameters. This extra level of parameters provide extra flexibility to a model, but generally requires more data to reach a well-posed problem. HBM refers to model classes with such a hierarchy (multiple levels) of parameters. Here, we define two types of HBM commonly found in the literature: a Hierarchical Prior Model (HPM) that further parameterizes the prior, and a Hierarchical Stochastic Model (HSM) that further parameterizes the stochastic model (or known as the likelihood function when evaluated at a given data).

HPM is a well-studied subject in the machine learning and compressive sensing community in the context of Sparse Bayesian Learning (SBL). Its applications include imaging [23, 6], gene selection [1], Bayesian compressive sensing [18, 17], etc. It is an effective tool to solve ill-posed regression problems by adding extra constraints in the prior distribution under the Bayesian framework. The well-known Bayesian Optimization method can also be interpreted using this model structure. On the other hand, the origin of HSM can be traced back to the multilevel model in the statistics community. One classic application is to analyze test results collected from multiple schools [11]. [24] and [16] are early publications on Bayesian analysis for multilevel regression, and Congdon [8] summarizes recent developments on this topic. Since most of the computational demand in Bayesian analysis comes from evaluating the stochastic model, HSM has a significantly larger computational cost than HPM. Although recently, there is increasing amount of HSM applications, the diversity of HSM uses is not comparable to HPM. However, recent developments in parallel computing opens up new opportunities to apply HSM to more complicated problems. In this paper, we discuss important theoretical implications of HSM and demonstrate an efficient approximation scheme with practical applications.

Despite the frequent use of the terminology "HBM" in the literature, the distinction between what we called the HSM and the HPM is often omitted. The overlapped usage of "HBM" has even caused confusion. For example, Guha et al. [13] studied the HPM. Nevertheless, the authors used [12] as a reference, which the type of HBM mentioned in [12] is actually the HSM. Hence, in this paper, we begin with a comparison between HPM and HSM in Section 2. Then, we turn our focus back to HSM and study its theoretical implications based on polynomial regression in Section 3. To tackle the high computational cost for analyzing HSM in practice, we propose an efficient approximation based on the idea of Importance Sampling and Empirical Interpolation Method [3] in Section 4. Section 5 includes two realistic examples using our approximation method. Finally, we give some concluding remarks in Section 6.

2 Classification of HBM: HPM versus HSM

For the ease of illustration, we begin with defining the notations that will be repeatedly used in this paper. To predict a quantity of interest y , one may have a stochastic forward model $F(x, \vec{\theta}, \vec{\epsilon})$ with model parameters $\vec{\theta}$ and the stochasticity part $\vec{\epsilon}$. This function defines the distribution $p(y|x, \vec{\theta})$. One of the most common example of $F(\cdot)$ is:

$$y = f(x, \vec{\theta}) + \epsilon_y, \quad (1)$$

where ϵ_y is chosen to follow a zero mean and σ_y standard deviation Gaussian distribution, $N(0, \sigma_y)$, for computational convenience. Given a set of input and output data pairs $\mathcal{D} = \{(\hat{x}_i, \hat{y}_i) | i, \dots, N_D\}$, the posterior distribution of $\vec{\theta}$ is inferred by the Bayes' Theorem:

$$p(\vec{\theta}|\mathcal{D}) = \frac{p(\mathcal{D}|\vec{\theta})p(\vec{\theta})}{p(\mathcal{D})}, \quad (2)$$

where the likelihood $p(\mathcal{D}|\vec{\theta}) \triangleq p(\hat{y}_1, \dots, \hat{y}_{N_D} | \hat{x}_1, \dots, \hat{x}_{N_D}, \vec{\theta})$ is the probability of observing the data values \mathcal{D} from the predictive model, the prior $p(\vec{\theta})$ is the initial belief of the values of $\vec{\theta}$, and the evidence (or marginal likelihood) $p(\mathcal{D}) \triangleq p(\hat{y}_1, \dots, \hat{y}_{N_D} | \hat{x}_1, \dots, \hat{x}_{N_D})$ is a critical term used in model selection [5, 4]. Then, a robust posterior prediction for an unobserved input-output pair (x_0, y_0) can be obtained by:

$$p(y_0|x_0, \mathcal{D}) = \int p(y_0|x_0, \vec{\theta})p(\vec{\theta}|\mathcal{D}) d\vec{\theta}. \quad (3)$$

HBM adds an extra level of hyperparameters $\vec{\psi}$ to this classical Bayesian model. This new hierarchy can be added to either the prior (HPM) or the stochastic model/likelihood (HSM), which will result in completely different model classes suitable for different problems. The major distinction of the two models comes from the different structure of information dependencies between all the stochastic variables. The graph representations shown in Figure 1 are very useful to display and study such relations [19]. The arrows denote the directions of the forward models, which implies that Bayes' Theorem is needed for inference in the opposite direction. Any nodes without an incoming arrow requires a prior distribution.

For HPM, an extra node is added to the starting nodes in order to parameterize the prior. This operation does not change the information dependency between the existing parameters. The posterior joint distribution of $\vec{\theta}$ and $\vec{\psi}$ is:

$$p(\vec{\theta}, \vec{\psi}|\mathcal{D}) = \frac{p(\mathcal{D}|\vec{\theta})p(\vec{\theta}|\vec{\psi})p(\vec{\psi})}{p(\mathcal{D})}. \quad (4)$$

If we integrate out $\vec{\psi}$ in the analysis, this is equivalent to recovering the classical Bayesian setting from HPM by choosing $p(\vec{\theta}) = \int p(\vec{\theta}|\vec{\psi})p(\vec{\psi}) d\vec{\psi}$, because

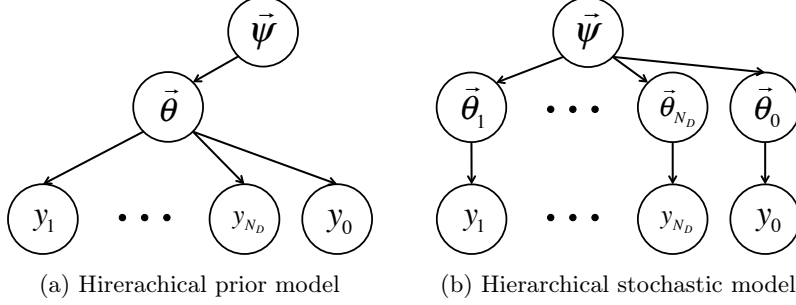


Figure 1: Graphical representations of the two HBM.

HPM does not affect the stochastic model part. Hence, the hyperparameters in HPM behave as latent variables. In most applications of HPM, we perform optimization of $\vec{\psi}$ by maximizing $p(\mathcal{D}|\vec{\psi})$, instead of the robust treatment of $\vec{\psi}$. For example, in the context of SBL, zero-mean Gaussian priors with different variances are chosen for the coefficients in a linear regression problem. This prior is called the Automatic Relevance Determination (ARD) prior that is proven to induce sparsity [20, 26]. A commonly used algorithm in SBL, called the Relevance Vector Machine, is based on this optimization setup [25, 27]. One can interpret this procedure as a Bayesian model selection problem in the continuous space.

For HSM, the hyperparameters affect the basic structure of the stochastic model. As shown in Figure 1b, each prediction is defined by a separate set of model parameters, and these parameters are correlated through the hyperparameters $\vec{\psi}$. Here, the original model parameters behave as latent variables, and $\vec{\psi}$ becomes the essential variables for future predictions y_0 . The posterior distribution of $\vec{\psi}$ is:

$$\begin{aligned}
 p(\vec{\psi}|\mathcal{D}) &= \frac{p(\mathcal{D}|\vec{\psi})p(\vec{\psi})}{p(\mathcal{D})} d\vec{\psi} \\
 \text{where } p(\mathcal{D}|\vec{\psi}) &= \prod_{i=1}^{N_D} p(D_i|\vec{\psi}), \\
 \text{and } p(D_i|\vec{\psi}) &= \int p(D_i|\vec{\theta}_i)p(\vec{\theta}_i|\vec{\psi}) d\vec{\theta}_i.
 \end{aligned} \tag{5}$$

Here, we denote the full data set as $\mathcal{D} = \{D_i|i = 1, \dots, N_D\}$, where D_i represents the observed data values for y_i . The choice of $p(\vec{\theta}_i|\vec{\psi})$ is flexible. For example, a statistical model is chosen in [12] and [8] and $\vec{\psi}$ is called the hyperparameter vector. The statistical model can be seen as a common prior for all $\vec{\theta}_i$. This explicit modeling for the prior separates different uncertainties by isolating out the uncertainty of $\vec{\theta}_i$ across multiple groups of predictions. It can give better predictions and also be related to causal inference [11]. Another example of such a hierarchical structure is found in a state-space model, where $\vec{\theta}_i$ is the state

variable and $p(\vec{\theta}_i|\vec{\psi})$ represents the underlying stochastic process. Therefore, $p(\vec{\theta}_i|\vec{\psi})$ can also be a theoretically informed model, instead of a purely empirical or statistical model.

In the current literatures, HPM is mainly used as a tool for optimal prior selection, while HSM is for explicit modeling of the uncertainty of model parameters across multiple groups of predictions for the data. However, HSM has the potential for more sophisticated modeling depending on the choice of $p(\vec{\theta}_i|\vec{\psi})$, though it comes with a significantly larger computational cost due to the extra integral shown in Equation 5. We note that since HPM and HSM are two exclusive types of HBM, in theory, they can co-exist in a complex HBM and may have more than one level of hyperparameters. It is the limit of computational power that constraints the usage of HBM to be only one level of hyperparameters. In this paper, we study the important theoretical implications of HSM and present an efficient and flexible approximation method for performing Bayesian analysis with HSM. We note that a similar comparison can be found in [22], but they present it purely from the perspective of a likelihood function of $\vec{\psi}$, instead of the HBM.

3 Theoretical Implications of HSM

Different assumptions on the uncertainties in a system can be represented using different hierarchical and non-hierarchical models. In this section, we verify that the effect of Occam’s Razor in the Bayesian model selection framework is applicable to both types of models. Moreover, we study the theoretical implications of using a HSM in three important aspects:

1. Separation of different types of uncertainties in a system (Section 3.1)
2. Identification of correlation between predictions or data (Section 3.1.2)
3. Uncertainty quantification of reduced order models (Section 3.2)

For computational accuracy, we demonstrate our results using polynomial regression because of the availability of many analytical solutions during the Bayesian analyses.

3.1 Separating different types of uncertainties using HSM

In this study, we demonstrate how HSM can capture different types of uncertainties using models that are embedded with different sources of stochasticity. We generate synthetic data based on two uncertainty models: (1) zero-mean Gaussian error ϵ_y added to the function $f(x, \vec{\theta})$, which can represent measurement noise in practice; (2) zero-mean Gaussian error ϵ_θ added to the model parameters $\vec{\theta}$, which can represent inherent model uncertainty due to, for example, insufficient knowledge of theoretically informed models or environmental variations across multiple experiments. We use a linear function $f(x, \theta) = \theta x$

(θ and x are scalars) and all data points are generated with independent ϵ_y and ϵ_θ . Three types of synthetic data are considered based on $\hat{\theta}$, a fixed value of the model parameter θ :

1. Additive error data, \mathcal{D}_1 : $y = f(x, \hat{\theta}) + \epsilon_y$, $\epsilon_y \sim N(\epsilon_y|0, \hat{\sigma}_y^2)$
2. Embedded error data, \mathcal{D}_{2a} : $y = f(x, \hat{\theta} + \epsilon_\theta)$, $\epsilon_\theta \sim N(\epsilon_\theta|0, \hat{\sigma}_\theta^2)$
3. Mixed error data, \mathcal{D}_{2b} : $y = f(x, \hat{\theta} + \epsilon_\theta) + \epsilon_y$, $\epsilon_y \sim N(\epsilon_y|0, \hat{\sigma}_y^2)$ and $\epsilon_\theta \sim N(\epsilon_\theta|0, \hat{\sigma}_\theta^2)$

where $N(z|\mu, \sigma^2)$ denotes a Gaussian distribution for z with mean μ and variance σ^2 . We set $\hat{\theta} = 1$, $\hat{\sigma}_\theta = 0.5$ and $\hat{\sigma}_y = 0.2$. Each data set contains 1000 data points independently generated from a uniformly distributed x value between preset bounds and random errors ϵ_y and ϵ_θ from the corresponding Gaussian distributions. We generate all three types of data twice, once with x between 0 and 1 and once with x between 0.4 and 1 (see Figures 2 and 3, respectively).

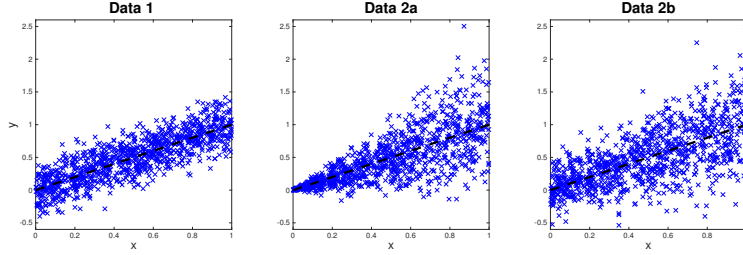


Figure 2: Three different types of contaminated data sets with x between 0 and 1. The black dash lines denote the actual function without any error.

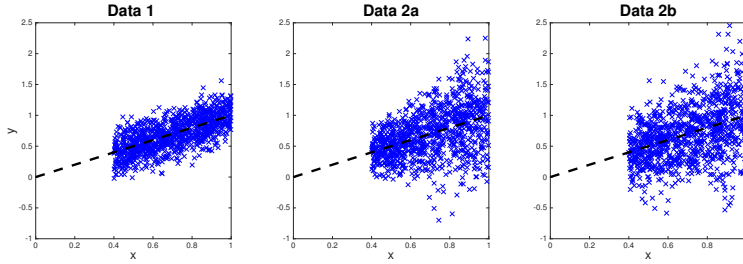


Figure 3: Three different types of contaminated data sets with x between 0.4 and 1. The black dash lines denote the actual function without any error.

Given the three types of data, we perform Bayesian model class selection on four different stochastic model classes:

1. Non-hierarchical model, \mathcal{M}_{1a} : This model class includes a uniform prior for θ between -1 and 3, and a Gaussian likelihood for θ , i.e., $p(\{x_i, y_i\}|\theta, \sigma_y) =$

$N(y_i|f(x_i, \theta), \sigma_y^2)$ for any i . Here, σ_y is also treated as a parameter to be inferred. We use a uniform prior for σ_y between 0.001 and 1.

2. HPM, \mathcal{M}_{1b} : This model class has the same likelihood as \mathcal{M}_{1a} and a hierarchical prior for θ . The prior is distributed as $N(\theta|\mu_\theta, \sigma_\theta^2)$, i.e., the hyperparameters $\vec{\psi} = \{\mu_\theta, \sigma_\theta\}$. We use uniform priors for μ_θ between -1 and 3 and σ_θ between 0.001 and 1. The prior for σ_y follows the one in \mathcal{M}_{1a} .
3. Zero noise HSM, \mathcal{M}_{2a} : This is an HSM with $p(\theta_i|\vec{\psi})$ modeled as the same Gaussian prior used for θ in \mathcal{M}_{1b} for all i . No additive error is assumed in this model. This implies that the likelihood function of θ_i is a delta function centered at the given data, i.e., $p(\{x_i, y_i\}|\theta_i) = \delta(y_i - f(x_i, \theta_i))$. The priors of the hyperparameters are the same as in \mathcal{M}_{1b} .
4. Full HSM, \mathcal{M}_{2b} : This is the same HSM as \mathcal{M}_{2a} except that a Gaussian additive error is assumed in this model. Hence, the likelihood function $p(\{x_i, y_i\}|\theta_i, \sigma_y) = N(y_i|f(x_i, \theta_i), \sigma_y^2)$. We use the same uniform priors as in \mathcal{M}_{1b} for σ_y , μ_θ and σ_θ .

Appendix A includes all derivations of the analytical expressions used for estimating the evidence $p(\mathcal{D}|\mathcal{M}_k)$, the marginal posterior probability density function (PDF) for the model parameters $p(\theta|\mathcal{D}, \mathcal{M}_k)$, and the posterior PDF for the hyperparameters $p(\psi|\mathcal{D}, \mathcal{M}_k)$, for a given model class \mathcal{M}_k . We use Monte Carlo Simulation to draw posterior samples of the hyperparameters for the "naive" estimation of the evidence. 10K samples are used to obtain accurate estimates with small variances.

3.1.1 Results and discussion

Table 1 and 2 summarize the results of the Bayesian inference and model selection for all six data sets (three with x between 0 and 1, and three with x between 0.4 and 1). The Bayesian model selection framework is capable of selecting the corresponding model used to generate the data. All models have a relatively good estimate for $\hat{\theta}$ given by $E[\theta|\mathcal{D}]$. The major difference appears in the accuracy of estimating $\hat{\sigma}_\theta$ and $\hat{\sigma}_y$ given by $Std[\theta|\mathcal{D}]$ and $E[\sigma_y|\mathcal{D}]$, respectively. In general, in terms of accurately estimating $\hat{\sigma}_\theta$ and $\hat{\sigma}_y$, \mathcal{M}_{1a} and \mathcal{M}_{1b} always prefer putting uncertainty into σ_y and thus perform well for \mathcal{D}_1 only. \mathcal{M}_{2a} cannot handle additive error and thus performs well for \mathcal{D}_{2a} only. \mathcal{M}_{2b} is the most flexible model and it performs relatively well for all data sets \mathcal{D}_1 , \mathcal{D}_{2a} and \mathcal{D}_{2b} .

We note that the results of \mathcal{M}_{1a} and \mathcal{M}_{1b} are extremely similar. This is because their only difference is the prior of θ , as explained in Section 2. In this case, the marginalized prior $p(\theta)$ in \mathcal{M}_{1b} is slightly larger than the prior in \mathcal{M}_{1a} . Hence, the evidence of both models are almost the same in all cases.

Furthermore, we observe that in the case of \mathcal{D}_{2b} with x between 0.4 and 1, \mathcal{M}_{2a} and \mathcal{M}_{2b} have a similar posterior model probability. Intuitively, the data

Table 1: Results for testing different stochastic model classes on the linear function (input x between 0 and 1). The row labeled “Ref.” shows the actual values for the corresponding estimates.

Data	Model	$E[\theta \mathcal{D}]$	$Std[\theta \mathcal{D}]$	$E[\sigma_y \mathcal{D}]$	$Std[\sigma_y \mathcal{D}]$	$\ln(\text{Evid.})$	$P(\mathcal{M}'_k \mathcal{D})$
\mathcal{D}_1	\mathcal{M}_{1a}	1.012	0.011	0.2030	0.0046	167.0	0.570
	\mathcal{M}_{1b}	1.012	0.011	0.2033	0.0049	166.7	0.424
	\mathcal{M}_{2a}	0.996	2.000	0.0000	0.0000	-2536.7	0.000
	\mathcal{M}_{2b}	1.017	0.133	0.1796	0.0008	162.4	0.006
	Ref.	1.000	0.000	0.200	- - -	- - -	- - -
\mathcal{D}_{2a}	\mathcal{M}_{1a}	1.002	0.016	0.2808	0.0065	-158.2	0.000
	\mathcal{M}_{1b}	1.002	0.016	0.2814	0.0067	-158.4	0.000
	\mathcal{M}_{2a}	0.978	0.525	0.0000	0.0000	285.3	1.000
	\mathcal{M}_{2b}	1.042	0.572	0.0021	0.0000	264.6	0.000
	Ref.	1.000	0.500	0.000	- - -	- - -	- - -
\mathcal{D}_{2b}	\mathcal{M}_{1a}	0.998	0.019	0.3334	0.0073	-328.6	0.000
	\mathcal{M}_{1b}	0.998	0.019	0.3330	0.0071	-328.7	0.000
	\mathcal{M}_{2a}	0.967	2.001	0.0000	0.0000	-16824.3	0.000
	\mathcal{M}_{2b}	1.027	0.453	0.1956	0.0041	-260.6	1.000
	Ref.	1.000	0.500	0.200	- - -	- - -	- - -

with small x values are more sensitive to additive noise because the noise-to-signal ratio is much higher than in the case of larger x values. When such data is not available, it is challenging to distinguish between \mathcal{D}_{2a} and \mathcal{D}_{2b} . This can be proved visually by observing the similarity between \mathcal{D}_{2a} and \mathcal{D}_{2b} in Figure 3 as compared to the one in Figure 2.

Our results suggest using \mathcal{M}_{2b} for model calibration, when no prior knowledge indicates that additive error is irrelevant. This HSM can appropriately separates the two types of uncertainties in our case study into the additive and non-additive parts based on the observed data. If computational power allows, Bayesian model class selection should always be performed among a set of different candidate models, as illustrated in this study.

3.1.2 Effect of grouping

In the previous section, all synthetic data is generated independently based on different stochastic models. In practice, we may group some of the predictions to be modeled under one set of parameter values $\vec{\theta}_i$, where i is the index for each group of predictions. For example, if an experiment is repeated in different laboratories, we may model the data from the same laboratory using the same set of model parameters, assuming the data shares the same environmental factors during the experiment. When the prediction grouping is not known, we may want to know the most plausible grouping, which defines a unique HSM.

We use the same linear function $f(x, \theta) = \theta x$ as before, but generate a new set of data \mathcal{D} . First of all, five random samples $\theta^{(i)}$ for $i = 1, \dots, 5$ are drawn from $N(\theta|\hat{\mu}_\theta, \hat{\sigma}_\theta^2)$. Then, 11 data points are generated for each $\theta^{(i)}$ to form a data set D_i by using the stochastic forward model $y = f(x, \theta^{(i)}) + \epsilon_y$, where ϵ_y is random error drawn from $N(\epsilon_y|0, \hat{\sigma}_y^2)$ independently for each data point. In

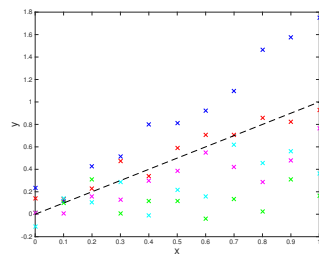
Table 2: Results for testing different stochastic model classes on the linear function (input x between 0.4 and 1). The row labeled “Ref.” shows the actual values for the corresponding estimates.

Data	Model	$E[\theta \mathcal{D}]$	$Std[\theta \mathcal{D}]$	$E[\sigma_y \mathcal{D}]$	$Std[\sigma_y \mathcal{D}]$	$\ln(\text{Evid.})$	$P(\mathcal{M}'_k \mathcal{D})$
\mathcal{D}_1	\mathcal{M}_{1a}	0.993	0.009	0.1985	0.0044	190.5	0.498
	\mathcal{M}_{1b}	0.993	0.009	0.1986	0.0050	190.5	0.486
	\mathcal{M}_{2a}	0.994	0.313	0.0000	0.0000	121.8	0.000
	\mathcal{M}_{2b}	0.984	0.035	0.2086	0.0017	187.0	0.015
	Ref.	1.000	0.000	0.200	- - -	- - -	- - -
\mathcal{D}_{2a}	\mathcal{M}_{1a}	0.998	0.016	0.3659	0.0082	-421.1	0.000
	\mathcal{M}_{1b}	0.997	0.016	0.3666	0.0072	-421.0	0.000
	\mathcal{M}_{2a}	1.000	0.500	0.0000	0.0000	-342.0	0.986
	\mathcal{M}_{2b}	1.006	0.469	0.0141	0.0001	-346.2	0.014
	Ref.	1.000	0.500	0.000	- - -	- - -	- - -
\mathcal{D}_{2b}	\mathcal{M}_{1a}	0.980	0.018	0.4136	0.0095	-542.4	0.000
	\mathcal{M}_{1b}	0.980	0.018	0.4136	0.0096	-542.7	0.000
	\mathcal{M}_{2a}	0.984	0.596	0.0000	0.0000	-508.9	0.558
	\mathcal{M}_{2b}	0.946	0.500	0.2101	0.0182	-509.1	0.442
	Ref.	1.000	0.500	0.200	- - -	- - -	- - -

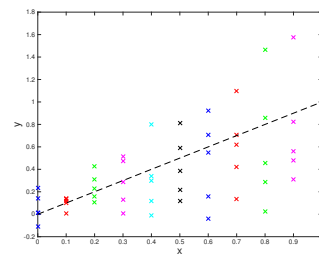
this study, $\hat{\mu}_\theta = 1$, $\hat{\sigma}_\theta = 0.5$ and $\hat{\sigma}_y = 0.1$. Figure 4a shows the five data sets used in the section. The sample mean and standard deviation of $\theta^{(i)}$ are 1.0191 and 0.4079, respectively.

To study the effect of grouping predictions in the HSM, we set up six different HSMs based on different information dependency between the predictions. We perform Bayesian model class selection to find the most plausible HSM to describe the data. The evidence is calculated similarly to \mathcal{M}_{2b} in Section 3.1. Details of the calculation can be found in Appendix A. The six data groupings are:

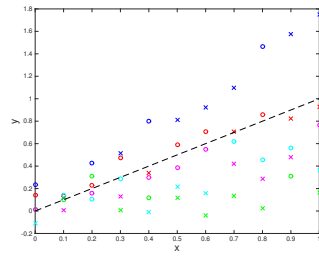
1. Actual grouping, \mathcal{M}'_1 : group predictions according to Figure 4a.
2. Constant x grouping, \mathcal{M}'_2 : predictions with the same x value are grouped together. This case represents grouping data according to the wrong variable. Figure 4b shows the resulting grouping.
3. 1/2 error grouping, \mathcal{M}'_3 : starting from \mathcal{M}'_1 , the predictions for each data set are further divided into 2 sets by the center line of the data points in the set. This center line is defined by the mean of θ inferred from each data point in the set. Figure 4c shows the resulting data grouping.
4. 1/4 error grouping, \mathcal{M}'_4 : starting from \mathcal{M}'_3 , the predictions for each data set are again further divided into 2 sets by the center line of the data points in the set, as done in \mathcal{M}'_3 .
5. Single prediction, \mathcal{M}'_5 : each prediction is treated as an independent data set (same as Section 3.1.1).
6. Random grouping, \mathcal{M}'_6 : predictions for the data points are randomly collected into five groups. Figure 4d shows the resulting grouping.



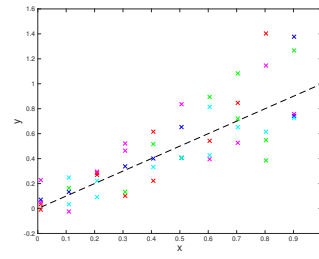
(a) Actual grouping



(b) Grouping across x



(c) Over grouping



(d) Random grouping

Figure 4: Different groupings of data sets. Different colors and marker combinations represent different data sets. The dashed lines denote the actual function without any error.

Table 3: Results for testing different HSM.

	$E[\theta \mathcal{D}]$	$Std[\theta \mathcal{D}]$	$E[\sigma_y \mathcal{D}]$	$Std[\sigma_y \mathcal{D}]$	$\ln(\text{Evid.})$	$P(\mathcal{M}'_k \mathcal{D})$
\mathcal{M}'_1	1.050	0.528	0.091	0.009	36.18	0.0012
\mathcal{M}'_2	1.051	0.061	0.230	0.023	-3.42	0.0000
\mathcal{M}'_3	1.036	0.473	0.066	0.007	42.88	0.9986
\mathcal{M}'_4	1.024	0.432	0.061	0.008	33.91	0.0001
\mathcal{M}'_5	1.052	0.348	0.097	0.020	5.49	0.0000
\mathcal{M}'_6	1.053	0.112	0.228	0.023	-3.14	0.0000
	$(\hat{\mu}_\theta)$	$(\hat{\sigma}_\theta)$	$(\hat{\sigma}_y)$			
<i>Ref.</i>	1.000	0.500	0.100	- - -	- - -	- - -

Table 3 summarizes the results of the Bayesian inference and model class selection for all six HSM. All models have a good estimate for $\hat{\mu}_\theta$ (corresponds to $E[\theta|\mathcal{D}]$). However, \mathcal{M}'_2 and \mathcal{M}'_6 , representing two completely wrong groupings, have significantly lower estimates for $\hat{\sigma}_\theta$ (given by $Std[\theta|\mathcal{D}]$), but larger estimates for the mean of $\hat{\sigma}_y$ (given by $E[\sigma_y|\mathcal{D}]$). Hence, they have a very low log-evidence value. \mathcal{M}'_1 has the closest estimates of $\hat{\mu}_\theta$, $\hat{\sigma}_\theta$ and $\hat{\sigma}_y$, but it is not the most probable model class among the six. \mathcal{M}'_3 is the most probable model class even though it has a low value of the estimates for both $\hat{\sigma}_\theta$ and $\hat{\sigma}_y$. We note that when the total number of data points is fixed, the more groups there are, the smaller the average number of predictions in a group is. This affects the likelihood of $\vec{\psi}$, which equals the product of $p(D_i|\vec{\psi})$ for each data set D_i (see Equation 5). When the number of $p(D_i|\vec{\psi})$ increases, the value of each evidence term may decrease. This decrease is due to a less peaked $p(D_i|\vec{\theta}_i)$ as the number of data points in the set decreases. The final evidence $p(\mathcal{D}|\mathcal{M}'_k)$ for a given data grouping \mathcal{M}'_k is a tradeoff between these two factors. Starting from \mathcal{M}'_1 being modified to \mathcal{M}'_3 , if we repeat the process many times, eventually we will reach \mathcal{M}'_5 . Therefore, \mathcal{M}'_1 , \mathcal{M}'_3 , \mathcal{M}'_4 and \mathcal{M}'_5 can be considered as a sequence of similar data grouping methods. In the end, the tradeoff suggests that \mathcal{M}'_3 is the most probable grouping. This implies that model class selection based on the evidence may not necessarily lead to the actual grouping, if it exists. Overall, it tries to minimize uncertainties in all parameters. On the other hand, if we choose the actual grouping, the uncertainty quantification will be accurate.

3.2 Uncertainty quantification for reduced order models using HSM

A reduced order model simplifies the information obtained in the original model. This loss of information can be treated as a source of uncertainty in the Bayesian framework. The structure of this type of uncertainty can be defined if the mapping between the original model and the reduced order model is completely known. In most cases, however, the mapping is not well-defined or too complicated to be studied directly. Therefore, there is no clear way to model such kind of loss of information. In this section, we compare the performance of using different model classes in Section 3.1 to represent the reduced order model

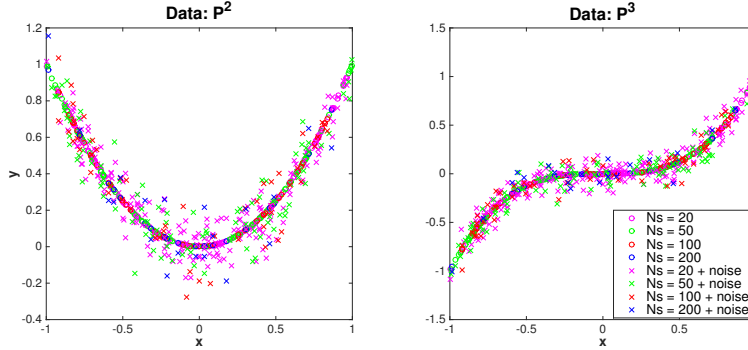


Figure 5: 8 data sets from $y = x^2$ (P^2) and 8 data sets from $y = x^3$ (P^3). The noises are additive Gaussian error with standard deviation $\hat{\sigma}_y = 0.1$.

uncertainty. This study is based on fitting second and third degree polynomial data with linear functions.

3.2.1 Problem setup for polynomials

We consider a total of 16 sets of data: 8 sets from a quadratic function $f(x) = x^2$ and 8 sets from a cubic function $f(x) = x^3$. For each type of function, we generate 20, 50, 100 and 200 data points twice (with and without noise), i.e., a total of 8 data sets. The x values of the data points are generated randomly in the interval $[-1, 1]$. Additive Gaussian error with standard deviation $\hat{\sigma}_y = 0.1$ is chosen for the noise. Figure 5 shows all 16 data sets used in this study.

We perform Bayesian model class selection and posterior robust prediction using three models described in Section 3.1: \mathcal{M}_{1a} (denoted as \mathcal{M}_1 in this section), \mathcal{M}_{2a} and \mathcal{M}_{2b} . Appendix A includes all derivations for the analytical expressions used to estimate the log-evidence values and posterior robust predictions $p(y|\mathcal{D}, \mathcal{M}_k)$ for a given model \mathcal{M}_k . We evaluate the posterior probability of observing each grid point (\hat{x}, \hat{y}) on a 2D fine grid of (x, y) using the analytical expressions in order to construct the distribution of the posterior robust prediction.

3.2.2 Results and discussion

Table 4 summarizes the results of model class selection for this study. \mathcal{M}_{2b} is the most probable model class for data with additive noise and \mathcal{M}_{2a} is the most probable model class for data without any additive noise. \mathcal{M}_1 is the least significant model in most cases because the error caused by the reduced order model is very non-linear. We do not expect that the additive noise alone is sufficient to explain the error. In this study, the sample size (number of data points) does not play a significant role in the results.

Figure 6 shows one set of the prediction results (the results are not sensitive

Table 4: Model class selection results for studying uncertainty quantification of reduced order models.

Sample Size		Quadratic (no noise)		Quadratic (with noise)		Cubic (no noise)		Cubic (with noise)	
		ln(Evid.)	$P(\mathcal{M}_k \mathcal{D})$	ln(Evid.)	$P(\mathcal{M}_k \mathcal{D})$	ln(Evid.)	$P(\mathcal{M}_k \mathcal{D})$	ln(Evid.)	$P(\mathcal{M}_k \mathcal{D})$
20	\mathcal{M}_1	-10.1	0.0000	-11.4	0.0095	5.9	0.0000	0.4	0.0908
	\mathcal{M}_{2a}	8.6	0.9873	-15.4	0.0002	18.4	0.9773	-8.9	0.0000
	\mathcal{M}_{2b}	4.3	0.0127	-6.8	0.9903	14.7	0.0227	2.7	0.9092
50	\mathcal{M}_1	-31.2	0.0000	-33.4	0.0018	28.5	0.0035	12.1	0.9625
	\mathcal{M}_{2a}	-6.2	0.9920	-1014.3	0.0000	34.2	0.9927	-237.2	0.0000
	\mathcal{M}_{2b}	-11.0	0.0080	-27.1	0.9982	28.6	0.0037	8.9	0.0375
100	\mathcal{M}_1	-67.1	0.0000	-66.2	0.0000	35.9	0.0000	21.3	0.0020
	\mathcal{M}_{2a}	-8.6	0.9685	-61.4	0.0000	60.8	0.9989	-31.2	0.0000
	\mathcal{M}_{2b}	-12.0	0.0315	-34.1	1.0000	53.9	0.0011	27.5	0.9980
200	\mathcal{M}_1	-101.9	0.0000	-102.3	0.0000	117.9	0.0000	80.6	0.0000
	\mathcal{M}_{2a}	57.7	0.9864	-272.5	0.0000	194.9	1.0000	-559.6	0.0000
	\mathcal{M}_{2b}	53.4	0.0136	-10.3	1.0000	170.0	0.0000	92.7	1.0000

to the change of the number of sample size). We observe that for cases without additive noise, predictions from \mathcal{M}_{2a} and \mathcal{M}_{2b} are almost the same. Because \mathcal{M}_{2a} is a simpler model than \mathcal{M}_{2b} (less parameters), Bayesian model class selection prefers \mathcal{M}_{2a} and gives it a higher log-evidence value. For cases with additive noise, the uncertainty of prediction from \mathcal{M}_{2a} is significantly larger than \mathcal{M}_{2b} . This is because θ is very sensitive to noise for x values close to zero when a model does not have any additive error component. As a result, Bayesian model class selection prefers \mathcal{M}_{2b} for its ability to better fit the data on average. These results suggest that the HSM with additive error is preferred as the first test for uncertainty quantification of reduced order model. The flexibility of such a model class to handle two types of uncertainty simultaneously (additive error and embedded error in the model parameters) results in a higher chance of discovering the underlying uncertainty structure of a reduced order model.

4 Efficient Approximation of HSM

The posterior distribution $p(\vec{\psi}|\mathcal{D})$ plays an important role in the HSM. Equation 5 shows that the calculation of the likelihood $p(\mathcal{D}|\vec{\psi})$ involves evaluations of multiple integrals, which correspond to the evidences for each data set D_i conditional on the hyperparameters $\vec{\psi}$. This leads to an extremely large computational cost for sampling from $p(\vec{\psi}|\mathcal{D})$, as well as estimating the evidence of the model, $p(\mathcal{D})$. Current approaches include using conjugate pairs for analytical results [8], approximating the integrals with Laplace Asymptotic Approximation [29], or using some advanced Markov Chain Monte Carlo techniques [21]. These methods are still restrictive for either studying many complex systems or the efficiency is not very scalable to handle extra data sets.

In this section, we present an efficient approximation method based on a special use of Importance Sampling. Most current research efforts focus on

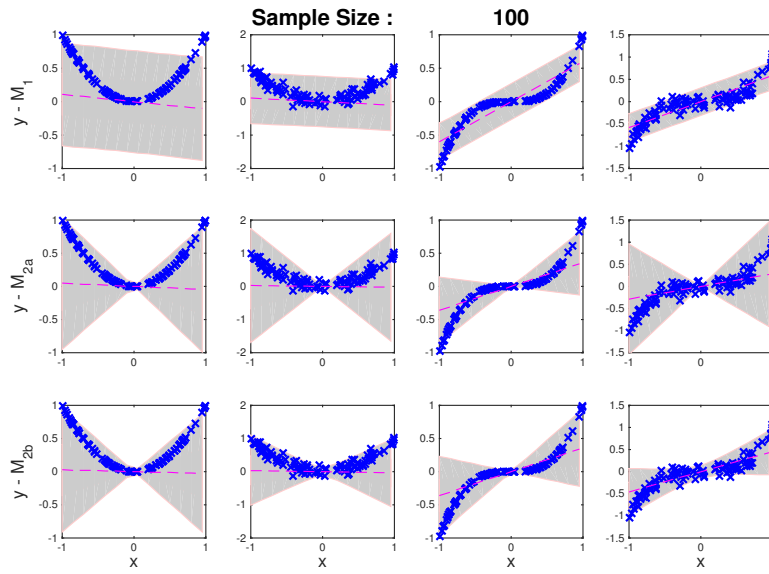


Figure 6: Posterior robust prediction for different models in different cases (sample size = 100). The purple dash lines denote the mean prediction and the grey area encloses 90% of the total probability density for the predicted value. Blue crosses are the data points. The four columns correspond to those in Table 4.

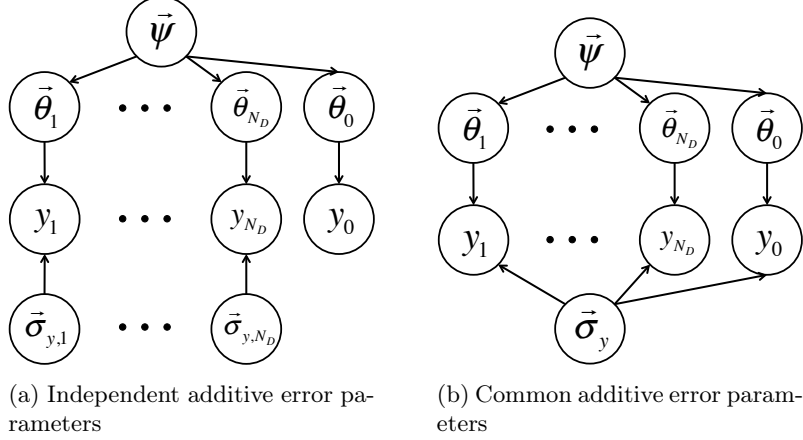


Figure 7: Graphical representations of the two alternative HSMs.

the HSM that adds only one extra level to the classical Bayesian model. Our method can be intuitively and efficiently extended to any complex HSMs with more levels of parameters. Moreover, we can lay out a standard procedure for fully analyzing any HSMs based on this method.

For ease of illustration, we demonstrate our method based on the commonly used stochastic model shown in Equation 1. In practice, there may be more complicated HSM than the one shown in Figure 1b. Here, we consider two possible alternatives (see Figure 7): (1) independent additive error parameters — applicable to heterogeneous data with different likelihoods [30], and (2) common additive error parameters — applicable to data sets that are expected to have the same measurement errors [29].

We introduce our method by first applying it to the basic HSM shown in Figure 1b, which represents the cases that either $\vec{\sigma}_y$ is known or it is uncertain and included with the other uncertain parameters $\vec{\theta}$. Then, we extend the idea to the two alternative HSMs shown in Figure 7.

4.1 Basic HSM

This model class considers σ_y to have a known value or to be treated as part of $\vec{\theta}$. Hence, we do not infer the distribution of σ_y from the data \mathcal{D} separately. Figure 1b shows the graphical representation for this case. In order to separate it from the two alternative HSMs in Figure 7, we denote this model class as

\mathcal{M}_{HS1} . The posterior distribution, the likelihood and the evidence for $\vec{\psi}$ are:

$$p(\psi|\mathcal{D}, \mathcal{M}_{HS1}) = \frac{p(\mathcal{D}|\psi, \mathcal{M}_{HS1})p(\psi|\mathcal{M}_{HS1})}{p(\mathcal{D}|\mathcal{M}_{HS1})} \quad (6)$$

$$\begin{aligned} p(\mathcal{D}|\psi, \mathcal{M}_{HS1}) &= \prod_{i=1}^{N_D} p(D_i|\psi, \mathcal{M}_{HS1}) \\ &= \prod_{i=1}^{N_D} \int p(D_i|\theta_i, \mathcal{M}_{HS1})p(\theta_i|\psi, \mathcal{M}_{HS1}) d\theta_i \end{aligned} \quad (7)$$

$$p(\mathcal{D}|\mathcal{M}_{HS1}) = \int p(\mathcal{D}|\psi, \mathcal{M}_{HS1})p(\psi|\mathcal{M}_{HS1}) d\psi \quad (8)$$

Sampling from $p(\psi|\mathcal{D}, \mathcal{M}_{HS1})$ requires repeated evaluations of the likelihood $p(\mathcal{D}|\psi, \mathcal{M}_{HS1})$ and thus evaluations of the integrals with different values of $\vec{\psi}$. Our idea is to approximate the integrals by Importance Sampling with a special choice of the proposal PDF that can significantly reduce the total computational cost. In many cases, the likelihood $p(D_i|\theta_i, \mathcal{M}_{HS1})$, which is related to the evaluation of $f(x, \vec{\theta})$, dominates the computational effort.

Our method begins by performing Bayesian inference for each data set D_i using the same likelihood $p(D_i|\theta_i, \mathcal{M}_{HS1})$, but a different prior (choice of such prior is discussed later). To be more specific, we draw samples $\{\vec{\theta}_i^{(j)}|j = 1, \dots, N_{s,i}\}$ from the posterior distribution $p(\theta_i|D_i, \mathcal{M}_i)$, where \mathcal{M}_i denotes this specific stochastic model class:

$$p(\vec{\theta}_i|D_i, \mathcal{M}_i) = \frac{p(D_i|\vec{\theta}_i, \mathcal{M}_i)p(\vec{\theta}_i|\mathcal{M}_i)}{p(D_i|\mathcal{M}_i)} \quad (9)$$

$$\text{where } p(D_i|\vec{\theta}_i, \mathcal{M}_i) = p(D_i|\vec{\theta}_i, \mathcal{M}_{HS1})$$

Then, we can approximate $p(D_i|\psi, \mathcal{M}_{HS1})$ based on IS with proposal distribution $q_i(\vec{\theta}_i) = p(\vec{\theta}_i|D_i, \mathcal{M}_i)$:

$$\begin{aligned} p(D_i|\psi, \mathcal{M}_{HS1}) &\approx \frac{1}{N_{s,i}} \sum_{j=1}^{N_{s,i}} \frac{p(D_i|\vec{\theta}_i^{(j)}, \mathcal{M}_{HS1})p(\vec{\theta}_i^{(j)}|\psi, \mathcal{M}_{HS1})}{q_i(\vec{\theta}_i^{(j)})} \\ &= \frac{p(D_i|\mathcal{M}_i)}{N_{s,i}} \sum_{j=1}^{N_{s,i}} \frac{p(\vec{\theta}_i^{(j)}|\psi, \mathcal{M}_{HS1})}{p(\vec{\theta}_i^{(j)}|\mathcal{M}_i)} \end{aligned} \quad (10)$$

where $\vec{\theta}_i^{(j)} \sim p(\vec{\theta}_i|D_i, \mathcal{M}_i)$

As a result, we only need to perform classical Bayesian inference once for each data set D_i (draw posterior samples and estimate the evidence $p(D_i|\mathcal{M}_i)$). Then, the hierarchical analysis comes as a post-processing of the results. A very good tool for such a Bayesian inference is the Transitional Markov Chain Monte Carlo (TMCMC) method, where the evidence comes as a by-product

of drawing posterior samples and the algorithm is inherently parallel [7]. The likelihood $p(\mathcal{D}|\psi, \mathcal{M}_{HS1})$ can then be estimated by:

$$p(\mathcal{D}|\vec{\psi}, \mathcal{M}_{HS1}) \approx \prod_{i=1}^{N_D} \left(\frac{p(D_i|M_i)}{N_{s,i}} \sum_{j=1}^{N_{s,i}} \frac{p(\vec{\theta}_i^{(j)}|\vec{\psi}, \mathcal{M}_{HS1})}{p(\vec{\theta}_i^{(j)}|\mathcal{M}_i)} \right), \text{ where } \vec{\theta}_i^{(j)} \sim p(\vec{\theta}_i|D_i, \mathcal{M}_i) \quad (11)$$

In other words, the evidence ratio between \mathcal{M}_{HS1} with given hyperparameters and \mathcal{M}_i can be approximated by the mean of the prior ratio between the two models over the posterior $\vec{\theta}_i$ samples of \mathcal{M}_i .

This approximation suffers the same problem as IS, i.e., the variance of the estimate depends strongly on the closeness of the integrand and the proposal distribution. In our case, they are exactly the same if $p(\vec{\theta}_i|\mathcal{M}_i) = p(\vec{\theta}_i|\vec{\psi}, \mathcal{M}_{HS1})$. Therefore, we should choose $p(\vec{\theta}_i|\mathcal{M}_i)$ to minimize its difference to $p(\vec{\theta}_i|\vec{\psi}, \mathcal{M}_{HS1})$ for all $\vec{\psi}$ weighted by the prior $p(\vec{\psi}|\mathcal{M}_{HS1})$. If we use the Kullback-Leibler divergence $D_{KL}(p||q)$ as a measure of difference between two distribution $p(x)$ and $q(x)$, this implies the choice:

$$p(\vec{\theta}_i|\mathcal{M}_i) = \operatorname{argmin} \int D_{KL}(p(\vec{\theta}_i|\vec{\psi}, \mathcal{M}_{HS1})||p(\vec{\theta}_i|\mathcal{M}_i))p(\vec{\psi}|\mathcal{M}_{HS1}) d\vec{\psi} \quad (12)$$

where $D_{KL}(p(x)||q(x)) = \int p(x) \log \frac{p(x)}{q(x)} dx$

A large number of posterior samples from $p(\vec{\theta}_i|D_i, \mathcal{M}_i)$ can also reduce the variance of the IS estimate. If this is computationally feasible, an alternative for $p(\theta_i|D_i, \mathcal{M}_i)$ would be a uniform distribution that covers the significant regions of $p(\vec{\theta}_i|\mathcal{M}_{HS1}) = \int p(\vec{\theta}_i|\vec{\psi}, \mathcal{M}_{HS1})p(\vec{\psi}|\mathcal{M}_{HS1}) d\vec{\psi}$. We note that the evidence will be insignificantly small for $\vec{\psi}$ values that lead to the prior $p(\vec{\theta}_i|\vec{\psi}, \mathcal{M}_{HS1})$ having a mode well away from the maximum of the likelihood $p(D_i|\vec{\theta}_i, \mathcal{M}_{HS1})$. Hence, accuracy is not important in this case. In the opposite case, IS is a good approximation. The inaccuracy problem is important only for the middle case, i.e., the prior mode is not too far and not too close to the maximum of the likelihood. The range of $\vec{\psi}$ leading to this important case depends on the posterior sample size $N_{s,i}$ for each data set D_i .

We note that this approximation method provides information about Bayesian inference for each individual data set first. This information is useful to compare with the HSM analysis to give further insight about the system of interest. Furthermore, our method is very efficient for introducing extra data sets because we do not need to rerun the whole problem. Instead, the overhead is only the classical Bayesian inferences for the extra data sets and a post-processing step that is not computationally intensive.

4.2 HSM with independent additive error parameters

This model class considers σ_y for the predictions of each data set to be independent of each other. Figure 7a shows the graphical representation for this case.

We denote this model class as \mathcal{M}_{HS2} . The posterior distribution, the likelihood and the evidence for $\vec{\psi}$ are:

$$p(\vec{\psi}|\mathcal{D}, \mathcal{M}_{HS2}) = \int p(\vec{\psi}, \vec{\sigma}_{y,1}, \dots, \vec{\sigma}_{y,N_D}|\mathcal{D}, \mathcal{M}_{HS2}) d\vec{\sigma}_y \dots d\vec{\sigma}_{y,N_D} \quad (13)$$

$$p(\vec{\psi}, \vec{\sigma}_{y,1}, \dots, \vec{\sigma}_{y,N_D}|\mathcal{D}, \mathcal{M}_{HS2}) = \frac{p(\mathcal{D}|\vec{\psi}, \vec{\sigma}_{y,1}, \dots, \vec{\sigma}_{y,N_D}, \mathcal{M}_{HS2})p(\vec{\psi}, \vec{\sigma}_{y,1}, \dots, \vec{\sigma}_{y,N_D}|\mathcal{M}_{HS2})}{p(\mathcal{D}|\mathcal{M}_{HS2})} \quad (14)$$

$$p(\mathcal{D}|\vec{\psi}, \vec{\sigma}_{y,1}, \dots, \vec{\sigma}_{y,N_D}, \mathcal{M}_{HS2}) = \prod_{i=1}^{N_D} \int p(D_i|\vec{\theta}_i, \vec{\sigma}_{y,i}, \mathcal{M}_{HS2})p(\vec{\theta}_i|\vec{\psi}, \mathcal{M}_{HS2}) d\vec{\theta}_i \quad (15)$$

$$\begin{aligned} & p(\mathcal{D}|\mathcal{M}_{HS2}) \\ &= \int p(\mathcal{D}|\vec{\psi}, \vec{\sigma}_{y,1}, \dots, \vec{\sigma}_{y,N_D}, \mathcal{M}_{HS2})p(\vec{\psi}, \vec{\sigma}_{y,1}, \dots, \vec{\sigma}_{y,N_D}|\mathcal{M}_{HS2}) d\vec{\psi} d\vec{\sigma}_{y,1} \dots d\vec{\sigma}_{y,N_D} \end{aligned} \quad (16)$$

The idea for \mathcal{M}_{HS1} does not apply to this model class directly because both the likelihood $p(D_i|\vec{\theta}_i, \vec{\sigma}_{y,i}, \mathcal{M}_{HS2})$ and the prior $p(\vec{\theta}_i|\vec{\psi}, \mathcal{M}_{HS2})$ will be affected by the sampling of $\vec{\psi}, \vec{\sigma}_{y,1}, \dots, \vec{\sigma}_{y,N_D}$. However, if prediction is the ultimate goal, Figure 7a implies that $\vec{\sigma}_{y,1}, \dots, \vec{\sigma}_{y,N_D}$ are not relevant as long as the posterior distribution of $\vec{\psi}$ is known. Then, we can again use IS to estimate $p(\vec{\psi}|\mathcal{D}, \mathcal{M}_{HS2})$ directly.

We note that this is a special case of \mathcal{M}_{HS1} with $p(\vec{\theta}_i, \vec{\sigma}_{y,i}|\vec{\psi}) = p(\vec{\theta}_i|\vec{\psi})p(\vec{\sigma}_{y,i})$, i.e., $\vec{\sigma}_{y,i}$ is independent of $\vec{\psi}$. Similar to \mathcal{M}_{HS1} , we first draw samples $\{(\vec{\theta}_i^{(j)}, \vec{\sigma}_{y,i}^{(j)})|j = 1, \dots, N_{s,i}\}$ from the posterior distribution $p(\vec{\theta}_i, \vec{\sigma}_{y,i}|D_i, \mathcal{M}_i)$:

$$p(\vec{\theta}_i, \vec{\sigma}_{y,i}|D_i, \mathcal{M}_i) = \frac{p(D_i|\vec{\theta}_i, \vec{\sigma}_{y,i}, \mathcal{M}_i)p(\vec{\theta}_i|\mathcal{M}_i)p(\vec{\sigma}_{y,i}|\mathcal{M}_i)}{p(D_i|\mathcal{M}_i)} \quad (17)$$

$$\text{where } p(D_i|\vec{\theta}_i, \vec{\sigma}_{y,i}, \mathcal{M}_i) = p(D_i|\vec{\theta}_i, \vec{\sigma}_{y,i}, \mathcal{M}_{HS2})$$

Following the procedure for \mathcal{M}_{HS1} , we can derive that:

$$\begin{aligned} p(D_i|\vec{\psi}, \mathcal{M}_{HS2}) &\approx \frac{p(D_i|\mathcal{M}_i)}{N_{s,i}} \sum_{j=1}^{N_{s,i}} \frac{p(\vec{\theta}_i^{(j)}|\vec{\psi}, \mathcal{M}_{HS2})}{p(\vec{\theta}_i^{(j)}|\mathcal{M}_i)} \frac{p(\vec{\sigma}_{y,i}^{(j)}|\mathcal{M}_{HS2})}{p(\vec{\sigma}_{y,i}^{(j)}|\mathcal{M}_i)} \\ &\text{where } (\vec{\theta}_i^{(j)}, \vec{\sigma}_{y,i}^{(j)}) \sim p(\vec{\theta}_i, \vec{\sigma}_{y,i}|D_i, \mathcal{M}_i) \end{aligned} \quad (18)$$

Again, the HSM analysis comes as a post-processing step and $p(\mathcal{D}|\vec{\psi}, \mathcal{M}_{HS2}) = \prod_{i=1}^{N_D} p(D_i|\vec{\psi}, \mathcal{M}_{HS2})$ can then be estimated using the posterior samples $(\vec{\theta}_i^{(j)}, \vec{\sigma}_{y,i}^{(j)})$:

$$p(D_i|\vec{\psi}, \mathcal{M}_{HS2}) \approx \prod_{i=1}^{N_D} \left(\frac{p(D_i|\mathcal{M}_i)}{N_{s,i}} \sum_{j=1}^{N_{s,i}} \frac{p(\vec{\theta}_i^{(j)}|\vec{\psi}, \mathcal{M}_{HS2})}{p(\vec{\theta}_i^{(j)}|\mathcal{M}_i)} \frac{p(\vec{\sigma}_{y,i}^{(j)}|\mathcal{M}_{HS2})}{p(\vec{\sigma}_{y,i}^{(j)}|\mathcal{M}_i)} \right) \quad (19)$$

As a result, we can directly obtain posterior samples from $p(\vec{\psi}|\mathcal{D}, \mathcal{M}_{HS2})$ using MCS or MCMC methods combined with Equation 19. We note that the estimation can be further simplified and more accurate when choosing $p(\sigma_{y,i}|\mathcal{M}_i) = p(\sigma_{y,i}|\mathcal{M}_{HS2})$ for all i .

4.3 HSM with common additive error parameters

This model class considers a single value of $\vec{\sigma}_y$ to be shared by all predictions. Figure 7b shows the graphical representation for this case. We denote this model class as \mathcal{M}_{HS3} . In this model class, $\vec{\sigma}_y$ affects future prediction as well. Hence, the approach for \mathcal{M}_{HS2} will not work. We adopt the idea from EIM to solve the problem. The basic concept of our approach is to approximate the likelihood function $p(D_i|\vec{\theta}_i, \vec{\sigma}_y, \mathcal{M}_{HS3})$ as a linear sum of multiple copies of the same term, but with fixed $\vec{\sigma}_y$ values. Then, we can perform the classical Bayesian analyses for those fixed $\vec{\sigma}_y$ cases and use the same IS approach as in \mathcal{M}_{HS1} and \mathcal{M}_{HS2} to approximate the HSM analysis for \mathcal{M}_{HS3} . First, we look at how to use the EIM idea to estimate the posterior distributions in \mathcal{M}_{HS3} . Then, we discuss how to train the linear approximations for the likelihood $p(D_i|\vec{\theta}_i, \vec{\sigma}_y, \mathcal{M}_{HS3})$.

We denote $g_l(\vec{\theta}_i) = p(D_i|\vec{\theta}_i, \vec{\sigma}_{y,l}, \mathcal{M}_{HS3})$ as a basis function for some fixed value of $\vec{\sigma}_{y,l}$ for $l = 1, \dots, L$. EIM assumes the following approximation for the fixed number of bases L and some coefficients $\alpha_l(\vec{\sigma}_y)$:

$$p(D_i|\vec{\theta}_i, \vec{\sigma}_y, \mathcal{M}_{HS3}) \approx \sum_{l=1}^L \alpha_l(\vec{\sigma}_y) g_l(\vec{\theta}_i) \quad (20)$$

Then, we can approximate the likelihood $p(D_i|\vec{\psi}, \vec{\sigma}_y, \mathcal{M}_{HS3}) = \int p(D_i|\vec{\theta}_i, \vec{\sigma}_y, \mathcal{M}_{HS3}) p(\vec{\theta}_i|\vec{\psi}, \mathcal{M}_{HS3}) d\vec{\theta}_i$ using Equation 20:

$$p(D_i|\vec{\psi}, \vec{\sigma}_y, \mathcal{M}_{HS3}) \approx \sum_{l=1}^L \alpha_l(\vec{\sigma}_y) \int p(D_i|\vec{\theta}_i, \vec{\sigma}_{y,l}, \mathcal{M}_{HS3}) p(\vec{\theta}_i|\vec{\psi}, \mathcal{M}_{HS3}) d\vec{\theta}_i \quad (21)$$

Since $\vec{\sigma}_{y,l}$ is fixed, we can apply the same IS approach as in \mathcal{M}_{HS2} and \mathcal{M}_{HS1} to the integrals in Equation 21. We draw samples $\{\vec{\theta}_i^{(j)}|j = 1, \dots, N_{s,i}\}$ from $p(\vec{\theta}_i|D_i, \vec{\sigma}_{y,l}, \mathcal{M}_i)$:

$$p(\vec{\theta}_i|D_i, \vec{\sigma}_{y,l}, \mathcal{M}_i) = \frac{p(D_i|\vec{\theta}_i, \vec{\sigma}_{y,l}, \mathcal{M}_i) p(\vec{\theta}_i|\mathcal{M}_i)}{p(D_i|\vec{\sigma}_{y,l}, \mathcal{M}_i)} \quad (22)$$

where $p(D_i|\vec{\theta}_i, \vec{\sigma}_{y,l}, \mathcal{M}_{HS3}) = p(D_i|\vec{\theta}_i, \vec{\sigma}_{y,l}, \mathcal{M}_i)$

Then, we use the proposal $q_{l,i}(\vec{\theta}_i) = p(\vec{\theta}_i|D_i, \vec{\sigma}_{y,l}, \mathcal{M}_i)$:

$$\begin{aligned}
p(D_i|\vec{\psi}, \vec{\sigma}_{y,l}, \mathcal{M}_{HS3}) &\approx \frac{1}{N_{s,i,l}} \sum_{j=1}^{N_{s,i,l}} \frac{p(D_i|\vec{\theta}_{i,l}^{(j)}, \vec{\sigma}_{y,l}, \mathcal{M}_{HS3}) p(\vec{\theta}_{i,l}^{(j)}|\vec{\psi}, \mathcal{M}_{HS3})}{q_{l,i}(\vec{\theta}_{i,l}^{(j)})} \\
&= \frac{p(D_i|\vec{\sigma}_{y,l}, \mathcal{M}_i)}{N_{s,i,l}} \sum_{j=1}^{N_{s,i,l}} \frac{p(\vec{\theta}_{i,l}^{(j)}|\vec{\psi}, \mathcal{M}_{HS3})}{p(\vec{\theta}_{i,l}^{(j)}|\mathcal{M}_i)} \\
&\text{where } \vec{\theta}_{i,l}^{(j)} \sim p(\vec{\theta}_i|D_i, \vec{\sigma}_{y,l}, \mathcal{M}_i)
\end{aligned} \tag{23}$$

As a result, the likelihood $p(\mathcal{D}|\vec{\psi}, \vec{\sigma}_y, \mathcal{M}_{HS3})$ can, then, be estimated by:

$$\begin{aligned}
p(\mathcal{D}|\vec{\psi}, \vec{\sigma}_y, \mathcal{M}_{HS3}) &\approx \prod_{i=1}^{N_D} \left(\sum_{l=1}^L \alpha_l(\vec{\sigma}_y) \frac{p(D_i|\vec{\sigma}_{y,l}, \mathcal{M}_i)}{N_{s,i,l}} \sum_{j=1}^{N_{s,i,l}} \frac{p(\vec{\theta}_{i,l}^{(j)}|\vec{\psi}, \mathcal{M}_{HS3})}{p(\vec{\theta}_{i,l}^{(j)}|\mathcal{M}_i)} \right) \\
&\text{where } \vec{\theta}_{i,l}^{(j)} \sim p(\vec{\theta}_i|D_i, \vec{\sigma}_{y,l}, \mathcal{M}_i)
\end{aligned} \tag{24}$$

With this analytical expression for the approximation, the remaining problem is how to pick $\vec{\sigma}_{y,l}$ and how to find $\alpha_l(\vec{\sigma}_y)$ when estimating the posterior distribution $p(\vec{\psi}, \vec{\sigma}_y|\mathcal{D}, \mathcal{M}_{HS3})$.

4.3.1 Training basis functions

Hesthaven et al. [15] suggest using an adaptive greedy algorithm to select the set of basis functions $g_l(\theta_i)$. Based on this idea, we develop an algorithm that simultaneously selects basis functions and collects posterior samples used in Equation 24. The greedy algorithm reduces the maximum absolute value of the error term $e_l(\vec{\theta}_i, \vec{\sigma}_{y,l})$ to some specified threshold \tilde{e}_{lim} over training sets of $\vec{\theta}_i$ and $\vec{\sigma}_y$, denote as Θ_i and Σ_y , respectively. The efficiency and accuracy of the algorithm on error estimation is a tradeoff that depends on the size of the training sets. The error term corresponding to L basis functions is defined as:

$$e_L(\vec{\theta}_i, \vec{\sigma}_y) = |p(D_i|\vec{\theta}_i, \vec{\sigma}_y, \mathcal{M}_{HS3}) - \sum_{l=1}^L \alpha_l(\vec{\sigma}_y) g_l(\vec{\theta}_i)| \tag{25}$$

Starting with some initial training sets $\Theta_i^{initial}$ and $\Sigma_y^{initial}$ and an initial set of basis functions $G^{initial}$, a new basis function with a corresponding $\vec{\sigma}_y$ value is chosen from Σ_y to maximize the error term e_L . Also, we record the $\vec{\theta}_i$ values that maximize e_L for the chosen $\vec{\sigma}_y$. These values will be used in the ‘‘online’’ estimation stage discussed later. Then, a set of posterior samples $\{\vec{\theta}_i\}_L$ is drawn using any MCMC methods with the chosen $\vec{\sigma}_y$, and the samples are added to the current training sets Θ_i . This process is repeated until the maximum error given Θ_i is below the threshold e_{lim} . In the end, we obtain a set of bases defined by $(\vec{\theta}_{i,l}, \vec{\sigma}_{y,l})$, $l = 1, \dots, L$, the posterior samples $\{\vec{\theta}_i\}_L$ and the evidence

values $p(D_i|\vec{\sigma}_{y,l}, \mathcal{M}_{HS3})$ for estimation of the approximate hierarchical Bayesian inference. Algorithm 1 summarizes this new adaptive training procedure.

Algorithm 1 : Adaptive greedy algorithm

Obtain initial $\vec{\sigma}_{y,1}, \dots, \vec{\sigma}_{y,L_0}$ for basis functions g_1, \dots, g_{L_0} from $G^{initial}$
Obtain initial $\vec{\theta}_{i,1}, \dots, \vec{\theta}_{i,L_0}$ for the corresponding basis functions from $G^{initial}$

Initialize counter $l \leftarrow L_0$
Initialize training sets $\Theta_i = \Theta_i^{initial}$ and $\Sigma_y = \Sigma_y^{initial}$
 $\tilde{\epsilon}_{L_0} = \max_{\vec{\sigma}_y \in \Sigma_y} \max_{\vec{\theta}_i \in \Theta_i} e_{L_0}(\vec{\theta}_i, \vec{\sigma}_y)$
while $\tilde{\epsilon}_l > \tilde{\epsilon}_{lim}$ **do**
 $l \leftarrow l + 1$
 $\vec{\sigma}_{y,l} = \operatorname{argmax}_{\vec{\sigma}_y \in \Sigma_y} \left\{ \max_{\vec{\theta}_i \in \Theta_i} e_{l-1}(\vec{\theta}_i, \vec{\sigma}_y) \right\}$
 $\vec{\theta}_{i,l} = \operatorname{argmax}_{\vec{\theta}_i \in \Theta_i} e_{l-1}(\vec{\theta}_i, \vec{\sigma}_{y,l})$
 $g_l(\vec{\theta}_i) = p(D_i|\vec{\theta}_i, \vec{\sigma}_{y,l}, \mathcal{M}_{HS3})$
 Obtain posterior samples $\{\vec{\theta}_i\}_l$ from distribution $p(\vec{\theta}_i|D_i, \vec{\sigma}_{y,l}, \mathcal{M}_i)$ using a MCMC method and calculate/record the evidence value $p(D_i|\vec{\sigma}_{y,l}, \mathcal{M}_i)$
 $\Theta_i \leftarrow \Theta_i \cup \{\vec{\theta}_i\}_l$
 $\tilde{\epsilon}_l = \max_{\vec{\sigma}_y \in \Sigma_y} \max_{\vec{\theta}_i \in \Theta_i} e_l(\vec{\theta}_i, \vec{\sigma}_y)$
end while

By enriching the training set Θ_i , the error estimate $e_L(\vec{\theta}_i, \vec{\sigma}_y)$ becomes more accurate and thus improves the EIM approximation. In many cases, most of the computational time for calculating $p(D_i|\vec{\theta}_i, \vec{\sigma}_y, \mathcal{M}_{HS3})$ is spent on evaluating the function $f(x, \vec{\theta}_i)$. Because this evaluation is already done during MCMC sampling, the overhead of using a larger training set Θ_i is relatively small. However, it is important to note that the improvement of the error estimate may saturate. This occurs when the different sets of posterior samples for different $\vec{\sigma}_y$ values all have the same high probability regions. In this case, there will be many redundant samples in the training set clustering around the peaks of $p(D_i|\vec{\theta}_i, \vec{\sigma}_y, \mathcal{M}_{HS3})$. To further improve the efficiency of this training algorithm, we can control the expansion of Θ_i such that we only add samples that are significantly different from the samples in the current Θ_i . A simple implementation is to monitor the spread of the chosen $\vec{\sigma}_y$ because similar $\vec{\sigma}_y$ values represents a similar likelihood function value $p(D_i|\vec{\theta}_i, \vec{\sigma}_y, \mathcal{M}_{HS3})$ and thus similar posterior samples are expected. Using the same argument, we suggest constructing the initial sets $\Theta_i^{initial}$ and $G^{initial}$ based on extreme values of $\vec{\sigma}_y$ (e.g., maximum and minimum values of $\vec{\sigma}_y$ in 1D case). Algorithm 2 constructs the initial sets based on this suggestion.

Algorithm 2 : Constructing initial sets $\Theta_i^{initial}, \Sigma_y^{initial}, G^{initial}$

Pick a fine grid for $\Sigma_y^{initial}$
Initialize $\vec{\sigma}_{y,1}, \dots, \vec{\sigma}_{y,K}$ as a sequence of extreme $\vec{\sigma}_y$ values ordered in ascending "peakness" of $p(D_i|\vec{\theta}_i, \vec{\sigma}_y, \mathcal{M}_{HS3})$, i.e., more probability content concentrated in a small region of $\vec{\theta}_i$ values (only qualitatively, need not to be accurate). This is the initial set of $\vec{\sigma}_y$ values that define the initial basis functions $g_1(\vec{\theta}_i), \dots, g_K(\vec{\theta}_i)$ in $G^{initial}$
Initialize $\Theta_i^{initial} \leftarrow$ empty set or a set of sparse grid points
for $k = 1$ **to** K **do**
 if $k = 1$ **then**
 $\vec{\theta}_{i,k} = \operatorname{argmax}_{\vec{\theta}_i \in \Theta_i} p(D_i|\vec{\theta}_i, \vec{\sigma}_{y,k}, \mathcal{M}_{HS3})$
 else
 $\vec{\theta}_{i,k} = \operatorname{argmax}_{\vec{\theta}_i \in \Theta_i} e_{l-1}(\vec{\theta}_i, \vec{\sigma}_{y,k})$
 end if
 Record $\vec{\theta}_{i,k}$ in $G^{initial}$ corresponding to basis function $g_k(\vec{\theta}_i)$
 Obtain posterior samples $\{\vec{\theta}_i\}_k$ from distribution $p(\vec{\theta}_i|D_i, \vec{\sigma}_{y,k}, \mathcal{M}_i)$ using a MCMC method and record the evidence value $p(D_i|\vec{\sigma}_{y,k}, \mathcal{M}_i)$
 $\Theta_i^{initial} \leftarrow \Theta_i^{initial} \cup \{\vec{\theta}_i\}_k$
end for

4.3.2 Online estimation

Once we obtain the set of bases $(\vec{\theta}_{i,l}, \vec{\sigma}_{y,l})$ that defines the basis functions $g_l(\vec{\theta}_i)$, where $l = 1, \dots, L$ and L is the total number of bases, we are ready to estimate $p(D_i|\vec{\psi}^{(s)}, \vec{\sigma}_y^{(s)}, \mathcal{M}_{HS3})$ with Equation 24 for any given hyperparameter sample $(\vec{\psi}^{(s)}, \vec{\sigma}_y^{(s)})$. First, we construct a matrix of the basis functions $g_{nl} = p(D_i|\vec{\theta}_{i,n}, \vec{\sigma}_{y,l}, \mathcal{M}_{HS3})$ for $n = 1, \dots, L$ and $l = 1, \dots, L$. Then, we calculate $P_n = p(D_i|\vec{\theta}_{i,n}, \vec{\sigma}_y^{(s)}, \mathcal{M}_{HS3})$, which is a vector of the actual values of the likelihood function. EIM constraints the linear approximation to be exact at the $\vec{\theta}_i$ bases. Hence, we find the vector of all $\alpha_l(\vec{\sigma}_y^{(s)})$ by solving the linear equations:

$$P_n = \sum_{l=1}^L \alpha_l(\vec{\sigma}_y^{(s)}) g_{nl}, \quad 1 \leq n \leq L \quad (26)$$

or in matrix form : $[g_{nl}]\{\alpha_l\} = \{P_n\}$

Once α_l is solved, we can include the recorded posterior samples $\{\vec{\theta}_i\}_l$ and evidence values $p(D_i|\vec{\sigma}_{y,l}, \mathcal{M}_i)$ corresponding to basis $g_l(\vec{\theta}_i)$ to perform fast estimation of $p(D_i|\vec{\psi}, \vec{\sigma}_y, \mathcal{M}_{HS3})$.

4.3.3 Numerical issue

In Equation 26, note that the vector $\{P_n\}$ and each column of the matrix $[g_{nl}]$ represents the likelihood value $p(D_i|\vec{\theta}_{i,n}, \vec{\sigma}_y, \mathcal{M}_{HS3})$ for different values of $\vec{\sigma}_y$.

The value of the Gaussian likelihood is very sensitive to $\vec{\sigma}_y$ and thus the values of the columns of $[g_{nl}]$ and $\{P_n\}$ may be of different orders of magnitude. As a result, the matrix inversion in Equation 26 often faces numerical issues because the matrix $[g_{nl}]$ is ill-conditioned. It is important to scale the matrix before solving the inversion problem. First, we rewrite the expression in Equation 26 as:

$$\sum_{l=1}^L (c_l^\alpha \hat{\alpha}_l) \cdot c_l^g \{\hat{g}_l\} = c^P \{\hat{P}_n\} \quad (27)$$

where $c^P \{\hat{P}_n\} = \{P_n\}$, $[g_{nl}] = [c_1^g \{\hat{g}_1\} \cdots c_L^g \{\hat{g}_L\}]$
and $\{\alpha_l\} = \{c_1^\alpha \hat{\alpha}_1 \cdots c_L^\alpha \hat{\alpha}_L\}^T$

If we choose $c_l^\alpha = c^P / c_l^g$ for all $l = 1, \dots, L$, the inversion problem becomes:

$$[\hat{g}_l] \{\hat{\alpha}_l\} = \{\hat{P}_n\} \quad (28)$$

We pick c^P and c_l^g to be the maximum value of the corresponding column vector. Then, the numerical issue due to scaling may not occur in Equation 28. We can first solve the inversion problem in Equation 28 and recover the coefficients α_l by $\{\alpha_l\} = \{c_1^\alpha \hat{\alpha}_1 \cdots c_L^\alpha \hat{\alpha}_L\}^T$ where $c_l^\alpha = c^P / c_l^g$. Also, it may be useful to store the values in log-scale during the entire computation process.

5 Illustrative Examples

We test our method using two examples: molecular dynamics simulation of Krypton and pharmacokinetic/pharmacodynamic (PKPD) modeling of a cancer drug. Both examples are taken from previous studies. Their results are used to validate the efficiency and accuracy of our method.

5.1 Molecular Dynamics: Krypton

Wu et al. [29] proposed using the HSM (called the HBM in the paper) to calibrate the Lennard-Jones (LJ) potential (ϵ_{LJ} and σ_{LJ}) of Krypton based on viscosity data \mathcal{D} reported from nine different laboratories. In that paper, the authors approximated the integrals in Equation 7 by Laplace Asymptotic Approximation (LAA). Since the joint posterior distribution of ϵ_{LJ} and σ_{LJ} for each data set is close to a Gaussian distribution, LAA is a good approximation for this particular case. We re-visit the problem with the same data and the same model class setup as described in [29]. Instead of the LAA approximation, we use our proposed method to estimate the joint posterior distribution of ϵ_{LJ} and σ_{LJ} and the posterior robust prediction.

Following the assumptions in [29], we use the HSM with common additive error. First, we build the EIM basis functions based on Section 4.3.1. We begin with a training set of ϵ_{LJ} and σ_{LJ} from a coarse grid of 256 points within the boundaries $100 \leq \epsilon_{LJ} \leq 400$ and $0.2 \leq \sigma_{LJ} \leq 0.5$. A fine grid of 172 points for σ_y is used between 0.0001 and 0.5, uniformly distributed in log scale. We

obtain a total of 50 basis functions for each data set to achieve an error less than 0.001%. Figure 8 shows the results of the EIM approximation. We observe that the bases of ϵ_{LJ} and σ_{LJ} coincide with the high likelihood values in the domain. This is consistent with our intuition that important regions of the likelihood should be included in the online estimation for better accuracy.

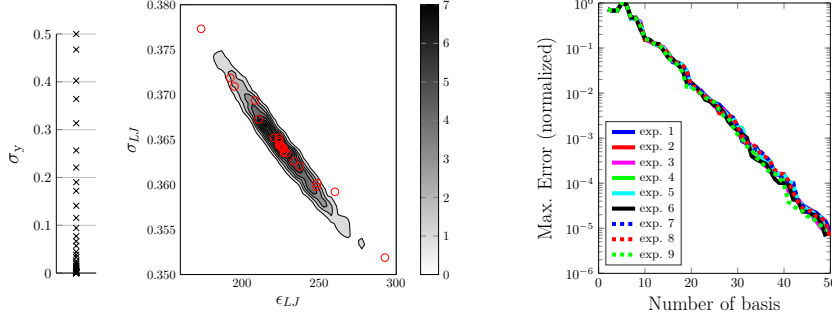


Figure 8: Results of the EIM approximation for the likelihood function. Left 2 plots: chosen bases of σ_y (black cross), ϵ_{LJ} and σ_{LJ} (red circles) for experiment 6. The gray scale contour shows the actual likelihood values for this experiment. Right plot: maximum error of the EIM estimate for each experiment as a function of the number of basis. The error is normalized by the maximum value of the actual function.

For each EIM basis, 2500 posterior samples of ϵ_{LJ} and σ_{LJ} are recorded to be used in the post-processing step of our HSM analysis. We use BASIS to draw the samples because it is highly parallel and the evidence is a by-product of the algorithm [31]. Then, we draw the posterior samples of the hyperparameters $\vec{\psi}$ with the post-processing step in our method. BASIS is once again used to draw 1000 posterior samples of $\vec{\psi}$. The posterior distribution of ϵ_{LJ} and σ_{LJ} is approximated by N posterior samples of $\vec{\psi}$:

$$\begin{aligned} p(\epsilon_{LJ}, \sigma_{LJ} | \mathcal{D}) &= \int p(\epsilon_{LJ}, \sigma_{LJ} | \vec{\psi}) p(\vec{\psi} | \mathcal{D}) d\vec{\psi} \\ &\approx \frac{1}{N} \sum_{i=1}^N p(\epsilon_{LJ}, \sigma_{LJ} | \vec{\psi}^{(i)}), \quad \text{where } \vec{\psi}^{(i)} \sim p(\vec{\psi} | \mathcal{D}) \end{aligned} \quad (29)$$

Figure 9 shows the results of the HSM analysis. We observe that our results are consistent with the results reported in [29]. This provides a verification of the accuracy of our method.

5.2 PKPD model: Cancer drug

Finley et al. [10] used the classical Bayesian model to study a PKPD model for the anti-vascular endothelial growth factor (anti-VEGF) cancer therapeutic agent, aflibercept. Multiple models with different number of parameters were

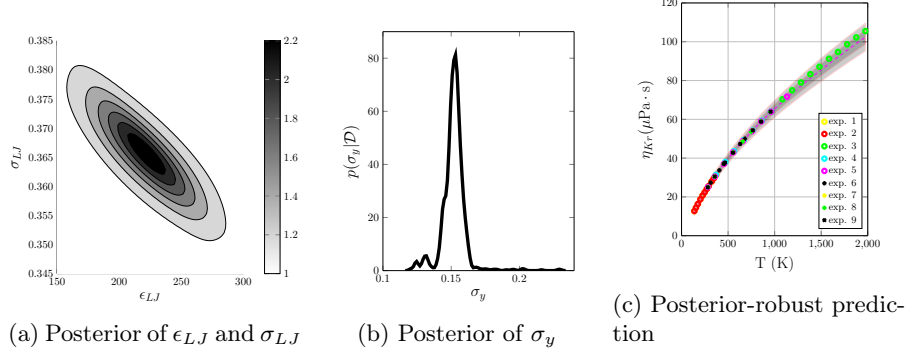


Figure 9: Posterior results for the Krypton MD simulation study.

calibrated using two types of clinical data: the plasma concentrations of free aflibercept and bound aflibercept. The data can be separated into 6 groups, each corresponding to a different dosage of the drug. We apply the HSM structure in Figure 1b to the basic model used in [10], which has three PK model parameters and five prediction error parameters. The prediction of the two types of data using these parameters are calculated by first running the system to reach steady state, and then imposing the corresponding drug dosages. The three model parameters $\{k_{EC}, k_N, k_T\}$ represent the secretion rate of VEGF in the blood, normal tissue and tumor compartment, respectively. The five prediction error parameters $\{\sigma_{SS}, \sigma_{BA_s}, \sigma_{BA_f}, \sigma_{FA_s}, \sigma_{FA_f}\}$ represent standard deviations of zero-mean Gaussian distributions for five different types of prediction errors: σ_{SS} — error on the steady state prediction, σ_{BA_s} — error on the bound aflibercept prediction scaled by time, σ_{BA_f} — error on the bound aflibercept prediction without any scaling, σ_{FA_s} — error on the free aflibercept prediction scaled by time, σ_{FA_f} — error on the free aflibercept prediction without any scaling. We choose independent Gaussian prior distributions for all of the eight parameters in log-10 scale.

Figure 10 and Table 5 compare the posterior distribution of the classical Bayesian model used in [10] with the posterior distribution of the basic HSM presented in this paper. We observe that fitting all data at once (the classical Bayesian model) leads to a significantly larger posterior variance for the model parameters, while the HSM leads to a larger posterior variance for the prediction error parameters. This implies that the knowledge about the model parameters is more transferable across different dosages than the knowledge about the prediction error parameters. We note that the steady state prediction is not affected by the change of dosage. Indeed, we observe that the inferred value of σ_{SS} in the HSM is similar to the value for the classical Bayesian model. The slight increase of the CV of σ_{SS} for the HSM case is expected because the data used for inference in the HSM is divided into groups of smaller data sets. Moreover, the scaled standard deviations σ_{BA_s} and σ_{FA_s} are larger, and the fixed standard deviations σ_{BA_f} and σ_{FA_f} are smaller in the HSM. This is also ex-

pected because in the classical Bayesian model, predictions for different dosages may have different level of sensitivity to the errors scaled by time. Hence, the classical Bayesian model tends to have a larger σ_{BAf} and σ_{FAf} than σ_{BA_s} and σ_{FA_s} . As a result, the posterior robust prediction based on the HSM (Figure 11) has completely different prediction errors for both bound affibercept and free affibercept compared to the one in [10].

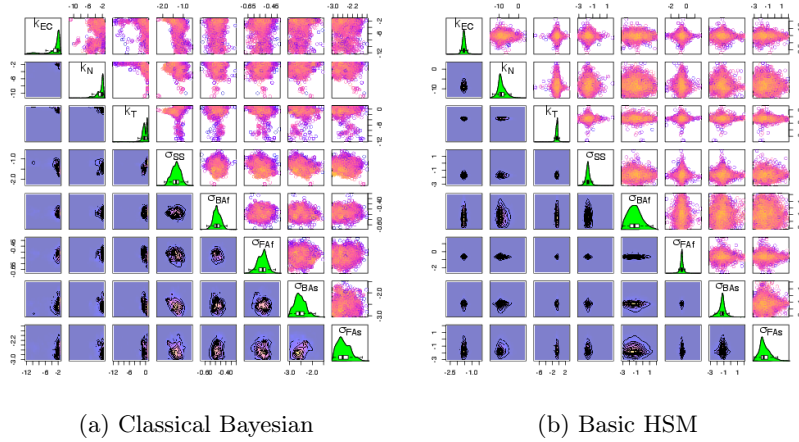


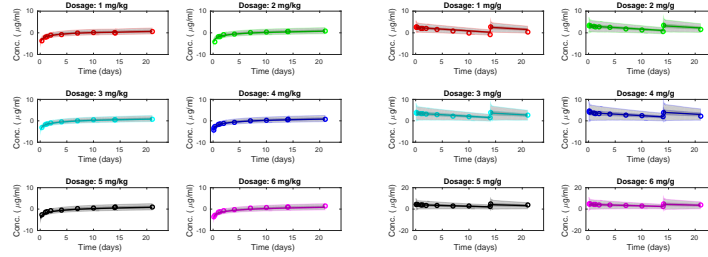
Figure 10: Posterior distributions for parameters of the basic PK model. Upper diagonal: projection of the posterior samples for all pairs of 2D parameter space (colors indicate log-likelihood values of the samples). Diagonal: marginal distributions of the model parameters estimated using kernel histograms. Box-plots denote the means and the 5 and 95 percentiles. Lower diagonal: projected densities in 2D parameter space constructed via a kernel estimate (coloring according to log-posterior values).

6 Conclusion

We demonstrate the benefits of using the hierarchical Bayesian framework in complex systems. Because of the double usage of the term HBM in the literature, we first explain the distinction between two very different HBMs: the HPM and the HSM. Then, we focus on studying the HSM, which has many interesting theoretical implications, as well as many computational challenges in practice. Based on examples of polynomial functions, we suggest that the HSM is capable of explicitly separating different types of uncertainties in a system, and can be an effective tool for modeling uncertainty of reduced order models. In order to apply HSM to practical problems, we propose an efficient approximation method to tackle the high computationally cost associated with using HSMs. Our method is a “bottom-up” approach that begins by drawing posterior samples of the model parameters for each data set. Then, we use a post-processing step to

Table 5: Posterior sample mean and coefficient of variation (CV) for parameters of the PK model. All samples are drawn in log-space. The mean is converted to the linear scale. The CV is calculated based on the samples in log-scale.

	Classical Bayesian		Basic HSM	
	Mean	CV(%)	Mean	CV(%)
k_{EC}	0.035	191.7	0.212	1.2
k_N	0.073	43.5	2×10^{-4}	52.0
k_T	0.634	442.7	0.535	23.6
σ_{SS}	0.257	3.4	0.183	7.5
σ_{BAf}	0.615	0.3	0.302	81.9
σ_{FAf}	0.590	0.3	0.534	11.0
σ_{BA_s}	0.080	2.4	0.276	26.2
σ_{FA_s}	0.077	2.7	0.235	44.2



(a) Bound aflibercept

(b) Free aflibercept

Figure 11: Posterior robust predictions for the plasma concentrations of aflibercept based on the basic HSM. Circles are the data points. Dark gray and light gray region are the 50% and 90% quantile range of the posterior distribution, respectively.

move up the hierarchy of the parameters for drawing posterior samples. Building on the basic HSM, we demonstrate our method using two alternative HSMs that have a more complicated hierarchical structure. Lastly, we validate our method using two illustrative examples based on previous studies: a molecular dynamics simulation of Krypton and a PKPD model of a cancer drug.

The HSM is a convenient and effective tool to build the stochastic model/likelihood for a complicated system. It also opens up new types of analyses and it results in different conclusions as compared to classical Bayesian inference. Our approximation method provides a standard procedure for analyzing hierarchical models. Beginning with an analysis for each individual data set, our method allows us to move up the hierarchy of the parameters efficiently to potentially extract more in-depth information about the system of interest. Moreover, it is very efficient for sequentially received data sets because our estimation scheme is a computationally fast post-processing step. So far, we have mainly used a

purely statistical (empirical) model for $p(\vec{\theta}|\vec{\psi})$. In future work, we plan to employ other choices for modeling $p(\vec{\theta}|\vec{\psi})$, which may motivate the development of new modeling tools.

A Derivations

We consider a set of data \mathcal{D} structured as $\mathcal{D} = \{D_i | i = 1, \dots, N_D\}$ with each subset of data $D_i = \{(x_{j,i}, y_{j,i}) | j = 1, \dots, N_{D_i}\}$, where $(x_{j,i}, y_{j,i})$ denotes a single data point. Using vector notations, we denote $\vec{x}_i = (x_{1,i}, \dots, x_{N_{D_i},i})^T$, $\vec{y}_i = (y_{1,i}, \dots, y_{N_{D_i},i})^T$, $\vec{x} = (\vec{x}_1^T, \dots, \vec{x}_{N_D}^T)^T$ and $\vec{y} = (\vec{y}_1^T, \dots, \vec{y}_{N_D}^T)^T$. Also, we denote the total number of data points $N_d = \sum_{i=1}^{N_D} N_{D_i}$. Hence, \vec{x} and \vec{y} are vectors with N_d elements.

Starting with a general formulation, we denote $\vec{\theta}$ as a vector of all model parameters and $\vec{\psi}$ as a vector of all hyperparameters. Repeatedly using Bayes' Theorem and the total probability theorem, we can derive the following expressions for a given model \mathcal{M}_k :

$$p(\vec{\psi}|\mathcal{D}, \mathcal{M}_k) = \frac{p(\mathcal{D}|\vec{\psi}, \mathcal{M}_k)p(\vec{\psi}|\mathcal{M}_k)}{p(\mathcal{D}|\mathcal{M}_k)} \quad (30)$$

$$p(\vec{\theta}|\mathcal{D}, \mathcal{M}_k) = \int p(\vec{\theta}|\mathcal{D}, \vec{\psi}, \mathcal{M}_k)p(\vec{\psi}|\mathcal{D}, \mathcal{M}_k) d\vec{\psi} \quad (31)$$

$$p(\vec{\theta}|\mathcal{D}, \vec{\psi}, \mathcal{M}_k) = \frac{p(\mathcal{D}|\vec{\theta}, \vec{\psi}, \mathcal{M}_k)p(\vec{\theta}|\vec{\psi}, \mathcal{M}_k)}{p(\mathcal{D}|\vec{\psi}, \mathcal{M}_k)} \quad (32)$$

$$p(\mathcal{M}_k|\mathcal{D}) = \frac{p(\mathcal{D}|\mathcal{M}_k)p(\mathcal{M}_k)}{p(\mathcal{D})} \quad (33)$$

$$p(\mathcal{D}|\mathcal{M}_k) = \int p(\mathcal{D}|\vec{\psi}, \mathcal{M}_k)p(\vec{\psi}|\mathcal{M}_k) d\vec{\psi} \quad (34)$$

$$p(\mathcal{D}|\vec{\psi}, \mathcal{M}_k) = \int p(\mathcal{D}|\vec{\theta}, \vec{\psi}, \mathcal{M}_k)p(\vec{\theta}|\vec{\psi}, \mathcal{M}_k) d\vec{\theta} \quad (35)$$

We note that usually $p(\mathcal{M}_k)$ is constant for all k because we do not want to introduce bias to any model before data is available. Hence, $p(\mathcal{M}_k|\mathcal{D}) \propto p(\mathcal{D}|\mathcal{M}_k)$. Also, $\vec{\theta}$ is simply a scalar θ in our study and the likelihood standard deviation σ_y is inferred separately. In the following, we derive the analytical expressions for the posterior distributions, the evidences and the robust-posterior predictions of the different models.

A.1 Non-hierarchical model, \mathcal{M}_{1a}

In this model, all data points are independent when θ and σ_y are known. By assuming a Gaussian distribution for the likelihood and a uniform prior $U(\theta) =$

$1/c_\theta$, we can derive that:

$$\begin{aligned} p(\mathcal{D}|\theta, \sigma_y, \mathcal{M}_{1a}) &= \prod_{i=1}^{N_D} \prod_{j=1}^{N_{D_i}} \frac{1}{\sqrt{2\pi}\sigma_y} \exp\left(-\frac{1}{2\sigma_y^2}(y_{j,i} - \theta x_{j,i})^2\right) \\ &= (2\pi)^{-N_d/2} \sigma_y^{-N_d} \exp\left(-\frac{1}{2\sigma_y^2}(\vec{y} - \theta \vec{x})^T(\vec{y} - \theta \vec{x})\right) \end{aligned} \quad (36)$$

Using Bayes' theorem, the posterior distribution is proportional to only the likelihood when the prior is uniform. It can be obtained from completing square for the exponential part of the Gaussian likelihood:

$$\begin{aligned} -\frac{1}{2\sigma_y^2}(\vec{y} - \theta \vec{x})^T(\vec{y} - \theta \vec{x}) &= -\frac{\vec{x}^T \vec{x}}{2\sigma_y^2} \left(\frac{\vec{y}^T \vec{y}}{\vec{x}^T \vec{x}} - 2\theta \frac{\vec{x}^T \vec{y}}{\vec{x}^T \vec{x}} + \theta^2 \right) \\ &= -\frac{1}{2\tilde{\sigma}_\theta^2} \left((\theta - \tilde{\mu}_\theta)^2 + \frac{\vec{y}^T \vec{y}}{\vec{x}^T \vec{x}} - \tilde{\mu}_\theta^2 \right) \quad (37) \\ \text{where } \tilde{\sigma}_\theta^2 &= \frac{\sigma_y^2}{\vec{x}^T \vec{x}}, \quad \tilde{\mu}_\theta = \frac{\vec{x}^T \vec{y}}{\vec{x}^T \vec{x}} = \frac{\vec{x}^T \vec{y}}{\sigma_y^2} \tilde{\sigma}_\theta^2 \end{aligned}$$

$$p(\mathcal{D}|\theta, \sigma_y, \mathcal{M}_{1a}) = \frac{1}{\sqrt{2\pi}\tilde{\sigma}_\theta} \exp\left(-\frac{(\theta - \tilde{\mu}_\theta)^2}{2\tilde{\sigma}_\theta^2}\right) \frac{\tilde{\sigma}_\theta}{(2\pi)^{\frac{N_d-1}{2}} \sigma_y^{N_d}} \exp\left(-\frac{\frac{\vec{y}^T \vec{y}}{\vec{x}^T \vec{x}} - \tilde{\mu}_\theta^2}{2\tilde{\sigma}_\theta^2}\right) \quad (38)$$

Hence, the posterior distribution is also Gaussian:

$$p(\theta|\mathcal{D}, \sigma_y, \mathcal{M}_{1a}) = N(\theta|\tilde{\mu}_\theta, \tilde{\sigma}_\theta^2), \quad \tilde{\sigma}_\theta = \frac{\sigma_y}{\sqrt{\vec{x}^T \vec{x}}}, \quad \tilde{\mu}_\theta = \frac{\vec{x}^T \vec{y}}{\vec{x}^T \vec{x}} = \frac{\vec{x}^T \vec{y}}{\sigma_y^2} \tilde{\sigma}_\theta^2 \quad (39)$$

and the evidence term $p(\mathcal{D}|\sigma_y, \mathcal{M}_{1a})$ calculated using Equation 35 (substitute $\vec{\psi}$ by σ_y) is:

$$p(\mathcal{D}|\sigma_y, \mathcal{M}_{1a}) = \frac{1}{c_\theta} \frac{\tilde{\sigma}_\theta}{(2\pi)^{\frac{N_d-1}{2}} \sigma_y^{N_d}} \exp\left(-\frac{1}{2} \left(\frac{\vec{y}^T \vec{y}}{\sigma_y^2} - \frac{\tilde{\mu}_\theta^2}{\tilde{\sigma}_\theta^2} \right)\right) \quad (40)$$

Then, we use N_s samples from Monte Carlo Simulation to estimate the posterior of σ_y and the model evidence:

$$p(\sigma_y|\mathcal{D}, \mathcal{M}_{1a}) \approx \sum_{j=1}^{N_s} w_j \delta(\sigma_y - \sigma_y^{(j)}) \quad (41)$$

$$\text{where } \sigma_y^{(j)} \sim p(\sigma_y|\mathcal{M}_{1a}), w_j \propto p(\mathcal{D}|\sigma_y^{(j)}, \mathcal{M}_{1a}), \sum_{j=1}^{N_s} w_j = 1$$

$$p(\mathcal{D}|\mathcal{M}_{1a}) \approx \frac{1}{N_s} \sum_{j=1}^{N_s} p(\mathcal{D}|\sigma_y^{(j)}, \mathcal{M}_{1a}) \quad (42)$$

The marginalized posterior of θ is estimated by substituting Equation 39 and 41 into Equation 31, which we substitute $\vec{\psi}$ by σ_y :

$$\begin{aligned}
p(\theta|\mathcal{D}, \mathcal{M}_{1a}) &\approx \sum_{j=1}^{N_s} w_j p(\theta|\mathcal{D}, \sigma_y^{(j)}, \mathcal{M}_{1a}) \\
&= \sum_{j=1}^{N_s} w_j N(\theta|\tilde{\mu}_\theta, (\tilde{\sigma}_\theta^{(j)})^2) \\
\text{where } \tilde{\mu}_\theta &= \frac{\vec{x}^T \vec{y}}{\vec{x}^T \vec{x}}, \quad \tilde{\sigma}_\theta^{(j)} = \frac{\sigma_y^{(j)}}{\sqrt{\vec{x}^T \vec{x}}}
\end{aligned} \tag{43}$$

As a result, the statistics of the marginalized posterior of θ and σ_y can also be estimated ($E[\cdot|\mathcal{D}]$ — posterior mean; $Std[\cdot|\mathcal{D}]$ — posterior standard deviation) based on the MCS samples:

$$\begin{aligned}
E[\theta|\mathcal{D}] &= \int \theta p(\theta|\mathcal{D}, \mathcal{M}_{1a}) d\theta \approx \tilde{\mu}_\theta \\
Std[\theta|\mathcal{D}] &\approx \sqrt{\left(\sum_{j=1}^{N_s} w_j \left((\tilde{\sigma}_\theta^{(j)})^2 + \tilde{\mu}_\theta^2 \right) \right) - E[\theta|\mathcal{D}]^2} \\
E[\sigma_y|\mathcal{D}] &\approx \sum_{j=1}^{N_s} w_j \sigma_y^{(j)} \\
Std[\sigma_y|\mathcal{D}] &= \sqrt{\left(\sum_{j=1}^{N_s} w_j (\sigma_y^{(j)})^2 \right) - E[\sigma_y|\mathcal{D}]^2}
\end{aligned} \tag{44}$$

To perform robust prediction of a new point (\hat{x}, \hat{y}) , we need to evaluate $p((\hat{x}, \hat{y})|\mathcal{D}, \mathcal{M}_{1a})$. Again, we can use the posterior samples from Equation 41:

$$\begin{aligned}
p((\hat{x}, \hat{y})|\mathcal{D}, \mathcal{M}_{1a}) &= \int p((\hat{x}, \hat{y})|\sigma_y, \mathcal{D}, \mathcal{M}_{1a}) p(\sigma_y|\mathcal{D}, \mathcal{M}_{1a}) d\sigma_y \\
&\approx \sum_{j=1}^{N_s} w_j p((\hat{x}, \hat{y})|\sigma_y^{(j)}, \mathcal{D}, \mathcal{M}_{1a})
\end{aligned} \tag{45}$$

where $p((\hat{x}, \hat{y})|\sigma_y^{(j)}, \mathcal{D}, \mathcal{M}_{1a})$ can be evaluated analytically based on the properties of the product of two Gaussian distributions (applied to the last line of this

equation):

$$\begin{aligned}
& p((\hat{x}, \hat{y}) | \sigma_y^{(j)}, \mathcal{D}, \mathcal{M}_{1a}) \\
&= \int p((\hat{x}, \hat{y}) | \theta, \sigma_y^{(j)}, \mathcal{M}_{1a}) p(\theta | \sigma_y^{(j)}, \mathcal{D}, \mathcal{M}_{1a}) d\theta \\
&= \int \frac{1}{\sqrt{2\pi}\sigma_y^{(j)}} \exp\left(-\frac{(\hat{y} - \theta\hat{x})^2}{2(\sigma_y^{(j)})^2}\right) \frac{1}{\sqrt{2\pi}\tilde{\sigma}_\theta^{(j)}} \exp\left(-\frac{(\theta - \tilde{\mu}_\theta)^2}{2(\tilde{\sigma}_\theta^{(j)})^2}\right) d\theta \\
&= \int \frac{1}{\hat{x}\sqrt{2\pi}(\sigma_y^{(j)}/\hat{x})} \exp\left(-\frac{(\theta - \hat{y}/\hat{x})^2}{2(\sigma_y^{(j)}/\hat{x})^2}\right) \frac{1}{\sqrt{2\pi}\tilde{\sigma}_\theta^{(j)}} \exp\left(-\frac{(\theta - \tilde{\mu}_\theta)^2}{2(\tilde{\sigma}_\theta^{(j)})^2}\right) d\theta \\
&= \frac{1}{\hat{x}\sqrt{2\pi}\left((\sigma_y^{(j)}/\hat{x})^2 + (\tilde{\sigma}_\theta^{(j)})^2\right)} \exp\left(-\frac{(\tilde{\mu}_\theta - \hat{y}/\hat{x})^2}{2\left((\sigma_y^{(j)}/\hat{x})^2 + (\tilde{\sigma}_\theta^{(j)})^2\right)}\right)
\end{aligned} \tag{46}$$

A.2 HPM, \mathcal{M}_{1b}

As explained in Section ??, this model is essentially the same as \mathcal{M}_{1a} , but with a different prior of θ . By assuming a Gaussian distribution for the likelihood and the prior of θ given the hyperparameters $\vec{\psi} = \{\mu_\theta, \sigma_\theta\}$, we can derive that:

$$\begin{aligned}
& p(\mathcal{D} | \sigma_y, \mathcal{M}_{1b}) p(\theta | \mu_\theta, \sigma_\theta, \mathcal{M}_{1b}) \\
&= \frac{1}{(2\pi)^{N_d/2} \sigma_y^{N_d}} \exp\left(-\frac{(\vec{y} - \theta\vec{x})^T (\vec{y} - \theta\vec{x})}{2\sigma_y^2}\right) \frac{1}{\sqrt{2\pi}\sigma_\theta} \exp\left(-\frac{(\theta - \mu_\theta)^2}{2\sigma_\theta^2}\right)
\end{aligned} \tag{47}$$

Using Bayes' theorem, the posterior distribution is proportional to the product of the likelihood and the prior. It can be obtained from completing square for the exponential part of the Gaussian distributions:

$$\begin{aligned}
& -\frac{1}{2\sigma_y^2} (\vec{y}^T \vec{y} - 2\theta \vec{x}^T \vec{y} + \theta^2 \vec{x}^T \vec{x}) - \frac{1}{2\sigma_\theta^2} (\theta^2 - 2\theta\mu_\theta + \mu_\theta^2) \\
&= -\frac{1}{2} \left(\theta^2 \left(\frac{\vec{x}^T \vec{x}}{\sigma_y^2} + \frac{1}{\sigma_\theta^2} \right) - 2\theta \left(\frac{\vec{x}^T \vec{y}}{\sigma_y^2} + \frac{\mu_\theta}{\sigma_\theta^2} \right) + \frac{\vec{y}^T \vec{y}}{\sigma_y^2} + \frac{\mu_\theta^2}{\sigma_\theta^2} \right) \\
&= -\frac{1}{2\tilde{\sigma}_\theta^2} \left(\theta^2 - 2\theta \left(\frac{\vec{x}^T \vec{y}\sigma_\theta^2 + \mu_\theta\sigma_y^2}{\vec{x}^T \vec{x}\sigma_\theta^2 + \sigma_y^2} \right) + \frac{\vec{y}^T \vec{y}\sigma_\theta^2 + \mu_\theta^2\sigma_y^2}{\vec{x}^T \vec{x}\sigma_\theta^2 + \sigma_y^2} \right) \\
&= -\frac{1}{2\tilde{\sigma}_\theta^2} \left((\theta - \tilde{\mu}_\theta)^2 + \frac{\vec{y}^T \vec{y}\sigma_\theta^2 + \mu_\theta^2\sigma_y^2}{\vec{x}^T \vec{x}\sigma_\theta^2 + \sigma_y^2} - \tilde{\mu}_\theta^2 \right)
\end{aligned} \tag{48}$$

$$\text{where } \tilde{\sigma}_\theta^2 = \frac{\sigma_\theta^2 \sigma_y^2}{\vec{x}^T \vec{x} \sigma_\theta^2 + \sigma_y^2}, \tilde{\mu}_\theta = \frac{\vec{x}^T \vec{y} \sigma_\theta^2 + \mu_\theta \sigma_y^2}{\vec{x}^T \vec{x} \sigma_\theta^2 + \sigma_y^2} = \left(\frac{\vec{x}^T \vec{y}}{\sigma_y^2} + \frac{\mu_\theta}{\sigma_\theta^2} \right) \tilde{\sigma}_\theta^2$$

$$\begin{aligned}
& p(\mathcal{D} | \sigma_y, \mathcal{M}_{1b}) p(\theta | \mu_\theta, \sigma_\theta, \mathcal{M}_{1b}) \\
&= \frac{1}{\sqrt{2\pi}\tilde{\sigma}_\theta} \exp\left(-\frac{(\theta - \tilde{\mu}_\theta)^2}{2\tilde{\sigma}_\theta^2}\right) \frac{\tilde{\sigma}_\theta}{(2\pi)^{\frac{N_d}{2}} \sigma_y^{N_d} \sigma_\theta} \exp\left(-\frac{1}{2} \left(\frac{\vec{y}^T \vec{y}}{\sigma_y^2} + \frac{\mu_\theta^2}{\sigma_\theta^2} - \frac{\tilde{\mu}_\theta^2}{\tilde{\sigma}_\theta^2} \right) \right)
\end{aligned} \tag{49}$$

Hence, the posterior distribution is also Gaussian:

$$p(\theta|\mathcal{D}, \mu_\theta, \sigma_\theta, \sigma_y, \mathcal{M}_{1b}) = N(\theta|\tilde{\mu}_\theta, \tilde{\sigma}_\theta^2)$$

$$\text{where } \tilde{\sigma}_\theta = \frac{\sigma_\theta \sigma_y}{\sqrt{\tilde{x}^T \tilde{x} \sigma_\theta^2 + \sigma_y^2}}, \tilde{\mu}_\theta = \frac{\tilde{x}^T \tilde{y} \sigma_\theta^2 + \mu_\theta \sigma_y^2}{\tilde{x}^T \tilde{x} \sigma_\theta^2 + \sigma_y^2} = \left(\frac{\tilde{x}^T \tilde{y}}{\sigma_y^2} + \frac{\mu_\theta}{\sigma_\theta^2} \right) \tilde{\sigma}_\theta^2 \quad (50)$$

and the evidence term $p(\mathcal{D}|\mu_\theta, \sigma_\theta, \sigma_y, \mathcal{M}_{1b})$ calculated using Equation 35 is:

$$p(\mathcal{D}|\mu_\theta, \sigma_\theta, \sigma_y, \mathcal{M}_{1b}) = \frac{\tilde{\sigma}_\theta}{(2\pi)^{\frac{N_d}{2}} \sigma_y^{N_d} \sigma_\theta} \exp \left(-\frac{1}{2} \left(\frac{\tilde{y}^T \tilde{y}}{\sigma_y^2} + \frac{\mu_\theta^2}{\sigma_\theta^2} - \frac{\tilde{\mu}_\theta^2}{\tilde{\sigma}_\theta^2} \right) \right) \quad (51)$$

We use N_s samples from Monte Carlo Simulation to estimate the posterior of μ_θ , σ_θ and σ_y :

$$p(\mu_\theta, \sigma_\theta, \sigma_y|\mathcal{D}, \mathcal{M}_{1b}) \propto p(\mathcal{D}|\mu_\theta, \sigma_\theta, \sigma_y, \mathcal{M}_{1b}) p(\mu_\theta, \sigma_\theta, \sigma_y|\mathcal{M}_{1b})$$

$$\approx \sum_{j=1}^{N_s} w_j \delta(\mu_\theta - \mu_\theta^{(j)}) \delta(\sigma_\theta - \sigma_\theta^{(j)}) \delta(\sigma_y - \sigma_y^{(j)}) \quad (52)$$

where $(\mu_\theta^{(j)}, \sigma_\theta^{(j)}, \sigma_y^{(j)}) \sim p(\mu_\theta, \sigma_\theta, \sigma_y|\mathcal{M}_{1b})$, $w_j \propto p(\mathcal{D}|\mu_\theta^{(j)}, \sigma_\theta^{(j)}, \sigma_y^{(j)}, \mathcal{M}_{1b})$, $\sum_{j=1}^{N_s} w_j = 1$

and the model evidence:

$$p(\mathcal{D}|\mathcal{M}_{1b}) \approx \frac{1}{N_s} \sum_{j=1}^{N_s} p(\mathcal{D}|\mu_\theta^{(j)}, \sigma_\theta^{(j)}, \sigma_y^{(j)}, \mathcal{M}_{1b}) \quad (53)$$

The marginalized posterior of θ can be estimated by substituting Equation 50 and 52 into Equation 31:

$$p(\theta|\mathcal{D}, \mathcal{M}_{1b}) \approx \sum_{j=1}^{N_s} w_j N(\theta|\tilde{\mu}_\theta^{(j)}, (\tilde{\sigma}_\theta^{(j)})^2)$$

$$\text{where } \tilde{\sigma}_\theta^{(j)} = \frac{\sigma_\theta^{(j)} \sigma_y^{(j)}}{\sqrt{\tilde{x}^T \tilde{x} (\sigma_\theta^{(j)})^2 + (\sigma_y^{(j)})^2}}, \tilde{\mu}_\theta^{(j)} = \frac{\tilde{x}^T \tilde{y} (\sigma_\theta^{(j)})^2 + \mu_\theta^{(j)} (\sigma_y^{(j)})^2}{\tilde{x}^T \tilde{x} (\sigma_\theta^{(j)})^2 + (\sigma_y^{(j)})^2} \quad (54)$$

As a result, the statistics of the marginalized posterior of θ and σ_y can also be estimated based on the MCS samples:

$$\begin{aligned}
E[\theta|\mathcal{D}] &= \int \theta p(\theta|\mathcal{D}, \mathcal{M}_{1b}) d\theta \approx \sum_{j=1}^{N_s} w_j \tilde{\mu}_\theta^{(j)} \\
Std[\theta|\mathcal{D}] &\approx \sqrt{\left(\sum_{j=1}^{N_s} w_j \left((\tilde{\sigma}_\theta^{(j)})^2 + (\tilde{\mu}_\theta^{(j)})^2 \right) \right) - E[\theta|\mathcal{D}]^2} \\
E[\sigma_y|\mathcal{D}] &\approx \sum_{j=1}^{N_s} w_j \sigma_y^{(j)} \\
Std[\sigma_y|\mathcal{D}] &= \sqrt{\left(\sum_{j=1}^{N_s} w_j (\sigma_y^{(j)})^2 \right) - E[\sigma_y|\mathcal{D}]^2}
\end{aligned} \tag{55}$$

To perform robust prediction of a new point (\hat{x}, \hat{y}) , we need to evaluate $p((\hat{x}, \hat{y})|\mathcal{D}, \mathcal{M}_{1b})$. Again, we can use the posterior samples from Equation 52:

$$\begin{aligned}
p((\hat{x}, \hat{y})|\mathcal{D}, \mathcal{M}_{1b}) &= \int p((\hat{x}, \hat{y})|\mu_\theta, \sigma_\theta, \sigma_y, \mathcal{D}, \mathcal{M}_{1b}) p(\mu_\theta, \sigma_\theta, \sigma_y|\mathcal{D}, \mathcal{M}_{1b}) d\mu_\theta d\sigma_\theta d\sigma_y \\
&\approx \sum_{j=1}^{N_s} w_j p((\hat{x}, \hat{y})|\mu_\theta^{(j)}, \sigma_\theta^{(j)}, \sigma_y^{(j)}, \mathcal{D}, \mathcal{M}_{1b})
\end{aligned} \tag{56}$$

where $p((\hat{x}, \hat{y})|\mu_\theta^{(j)}, \sigma_\theta^{(j)}, \sigma_y^{(j)}, \mathcal{D}, \mathcal{M}_{1b})$ can be evaluated analytically similar to the case in \mathcal{M}_{1a} :

$$\begin{aligned}
&p((\hat{x}, \hat{y})|\mu_\theta^{(j)}, \sigma_\theta^{(j)}, \sigma_y^{(j)}, \mathcal{D}, \mathcal{M}_{1b}) \\
&= \frac{1}{\sqrt{2\pi \left((\sigma_y^{(j)}/\hat{x})^2 + (\tilde{\sigma}_\theta^{(j)})^2 \right)}} \exp \left(-\frac{(\tilde{\mu}_\theta^{(j)} - \hat{y}/\hat{x})^2}{2 \left((\sigma_y^{(j)}/\hat{x})^2 + (\tilde{\sigma}_\theta^{(j)})^2 \right)} \right)
\end{aligned} \tag{57}$$

A.3 Zero noise HSFM, \mathcal{M}_{2a}

In this model, we assume the same Gaussian prior of θ_i for all $i = 1, \dots, N_D$, where θ_i corresponds to the data set $D_i \in \mathcal{D}$. The hyperparameters $\vec{\psi} = \{\mu_\theta, \sigma_\theta\}$ define the mean and the standard deviation of the Gaussian prior, respectively. By assuming a delta function for the likelihood, we can derive that for each data

set D_i :

$$\begin{aligned}
p(\theta_i|D_i, \mu_\theta, \sigma_\theta, \mathcal{M}_{2a}) &\propto p(D_i|\theta_i, \mathcal{M}_{2a})p(\theta_i|\mu_\theta, \sigma_\theta, \mathcal{M}_{2a}) \\
&= \prod_{j=1}^{N_{D_i}} \delta(y_{j,i} - \theta_i x_{j,i}) N(\theta_i|\mu_\theta, \sigma_\theta^2) \\
\text{i.e., } p(\theta_i|D_i, \mu_\theta, \sigma_\theta, \mathcal{M}_{2a}) &= \begin{cases} 1, & \text{if there exists } \theta_i \text{ such that } y_{j,i} = \theta_i x_{j,i} \text{ for all } j = 1, \dots, N_{D_i} \\ 0, & \text{otherwise} \end{cases}
\end{aligned} \tag{58}$$

And when the posterior equals 1, the evidence term $p(D_i|\mu_\theta, \sigma_\theta, \mathcal{M}_{2a})$ calculated using Equation 35 is:

$$p(D_i|\mu_\theta, \sigma_\theta, \mathcal{M}_{2a}) = \frac{1}{\sqrt{2\pi}\sigma_\theta} \exp\left(-\frac{1}{2\sigma_\theta^2} \left(\frac{y_{1,i}}{x_{1,i}} - \mu_\theta\right)^2\right) \cdot \prod_{j=1}^{N_{D_i}} \frac{1}{|x_{j,i}|} \tag{59}$$

Otherwise, the evidence is zero. Note that the extra product of inverse of $x_{j,i}$ comes from the Jacobian of the delta likelihood function.

Similar to \mathcal{M}_{1b} , we use N_s samples from Monte Carlo Simulation to estimate the posterior of μ_θ and σ_θ :

$$\begin{aligned}
p(\mu_\theta, \sigma_\theta|\mathcal{D}, \mathcal{M}_{2a}) &\approx \sum_{j=1}^{N_s} w_j \delta(\mu_\theta - \mu_\theta^{(j)}) \delta(\sigma_\theta - \sigma_\theta^{(j)}) \\
\text{where } (\mu_\theta^{(j)}, \sigma_\theta^{(j)}) &\sim p(\mu_\theta, \sigma_\theta|\mathcal{M}_{2a}), w_j \propto \prod_{i=1}^{N_D} p(D_i|\mu_\theta^{(j)}, \sigma_\theta^{(j)}, \mathcal{M}_{2a}), \sum_{j=1}^{N_s} w_j = 1
\end{aligned} \tag{60}$$

and the model evidence:

$$p(\mathcal{D}|\mathcal{M}_{2a}) \approx \frac{1}{N_s} \sum_{j=1}^{N_s} \left(\prod_{i=1}^{N_D} p(D_i|\mu_\theta^{(j)}, \sigma_\theta^{(j)}, \mathcal{M}_{2a}) \right) \tag{61}$$

In the HSFM, θ for future predictions (denoted as θ_0 in Section ??), is independent of the data when $\vec{\psi}$ is given. Hence, the marginalized posterior of θ depends only on the posterior of μ_θ and σ_θ and the prior of θ given μ_θ and σ_θ :

$$\begin{aligned}
p(\theta|\mathcal{D}, \mathcal{M}_{2a}) &= \int p(\theta|\mu_\theta, \sigma_\theta, \mathcal{M}_{2a}) p(\mu_\theta, \sigma_\theta|\mathcal{D}, \mathcal{M}_{2a}) d\mu_\theta d\sigma_\theta \\
&\approx \sum_{j=1}^{N_s} w_j N(\mu_\theta^{(j)}, \sigma_\theta^{(j)})
\end{aligned} \tag{62}$$

As a result, the statistics of the marginalized posterior of θ can also be estimated

based on the MCS samples:

$$\begin{aligned}
E[\theta|\mathcal{D}] &= \int \theta p(\theta|\mathcal{D}, \mathcal{M}_{2a}) d\theta \approx \sum_{j=1}^{N_s} w_j \mu_{\theta}^{(j)} \\
Std[\theta|\mathcal{D}] &\approx \sqrt{\left(\sum_{j=1}^{N_s} w_j \left((\sigma_{\theta}^{(j)})^2 + (\mu_{\theta}^{(j)})^2 \right) \right) - E[\theta|\mathcal{D}]^2}
\end{aligned} \tag{63}$$

To perform robust prediction of a new point (\hat{x}, \hat{y}) , we need to evaluate $p((\hat{x}, \hat{y})|\mathcal{D}, \mathcal{M}_{2a})$. Again, we can use the posterior samples from Equation 60:

$$\begin{aligned}
p((\hat{x}, \hat{y})|\mathcal{D}, \mathcal{M}_{2a}) &= \int p((\hat{x}, \hat{y})|\mu_{\theta}, \sigma_{\theta}, \mathcal{M}_{2a}) p(\mu_{\theta}, \sigma_{\theta}|\mathcal{D}, \mathcal{M}_{2a}) d\mu_{\theta} d\sigma_{\theta} \\
&\approx \sum_{j=1}^{N_s} w_j p((\hat{x}, \hat{y})|\mu_{\theta}^{(j)}, \sigma_{\theta}^{(j)}, \mathcal{M}_{2a})
\end{aligned} \tag{64}$$

Note that an unique feature of the HSFM is that future predictions do not depend on the data once all hyperparameters are given. Hence, $p((\hat{x}, \hat{y})|\mu_{\theta}^{(j)}, \sigma_{\theta}^{(j)}, \mathcal{M}_{2a})$ does not have \mathcal{D} in the conditional part, and we can evaluate it using Equation 59 by substituting D_i with (\hat{x}, \hat{y}) .

A.4 Full HSFM, \mathcal{M}_{2b}

In this model, we use the same setup as in \mathcal{M}_{2a} for the hyperparameters $\vec{\psi}$. By assuming a Gaussian distribution for the likelihood and the priors of θ_i for all $i = 1, \dots, N_D$ conditional on the hyperparameters, we can obtain analytical expressions similar to the one in \mathcal{M}_{1b} . One difference is that in this model, we now work with multiple data sets $D_i \in \mathcal{D}, i = 1, \dots, N_D$. Following a similar derivation as in \mathcal{M}_{1b} , the posterior distribution $p(\theta_i|D_i, \mu_{\theta}, \sigma_{\theta}, \sigma_y, \mathcal{M}_{2b})$ is Gaussian:

$$\begin{aligned}
p(\theta_i|D_i, \mu_{\theta}, \sigma_{\theta}, \sigma_y, \mathcal{M}_{2b}) &= N(\theta_i|\tilde{\mu}_{\theta,i}, (\tilde{\sigma}_{\theta,i})^2) \\
\text{where } \tilde{\sigma}_{\theta,i} &= \frac{\sigma_{\theta}\sigma_y}{\sqrt{\vec{x}_i^T \vec{x}_i \sigma_{\theta}^2 + \sigma_y^2}}, \tilde{\mu}_{\theta,i} = \frac{\vec{x}_i^T \vec{y}_i \sigma_{\theta}^2 + \mu_{\theta} \sigma_y^2}{\vec{x}_i^T \vec{x}_i \sigma_{\theta}^2 + \sigma_y^2} = \left(\frac{\vec{x}_i^T \vec{y}_i}{\sigma_y^2} + \frac{\mu_{\theta}}{\sigma_{\theta}^2} \right) \tilde{\sigma}_{\theta,i}^2
\end{aligned} \tag{65}$$

and the evidence term $p(D_i|\mu_{\theta}, \sigma_{\theta}, \sigma_y, \mathcal{M}_{2b})$ calculated using Equation 35 is:

$$p(D_i|\mu_{\theta}, \sigma_{\theta}, \sigma_y, \mathcal{M}_{2b}) = \frac{\tilde{\sigma}_{\theta,i}}{(2\pi)^{\frac{N_{D_i}}{2}} \sigma_y^{N_{D_i}} \sigma_{\theta}} \exp \left(-\frac{1}{2} \left(\frac{\vec{y}_i^T \vec{y}_i}{\sigma_y^2} + \frac{\mu_{\theta}^2}{\sigma_{\theta}^2} - \frac{\tilde{\mu}_{\theta,i}^2}{\tilde{\sigma}_{\theta,i}^2} \right) \right) \tag{66}$$

Similar to \mathcal{M}_{1b} , we use N_s samples from Monte Carlo Simulation to estimate the posterior of μ_θ , σ_θ and σ_y :

$$\begin{aligned}
p(\mu_\theta, \sigma_\theta, \sigma_y | \mathcal{D}, \mathcal{M}_{2b}) &\propto p(\mathcal{D} | \mu_\theta, \sigma_\theta, \sigma_y, \mathcal{M}_{2b}) p(\mu_\theta, \sigma_\theta, \sigma_y | \mathcal{M}_{2b}) \\
&= \prod_{i=1}^{N_D} p(D_i | \mu_\theta, \sigma_\theta, \sigma_y, \mathcal{M}_{2b}) p(\mu_\theta, \sigma_\theta, \sigma_y | \mathcal{M}_{2b}) \\
&\approx \sum_{j=1}^{N_s} w_j \delta(\mu_\theta - \mu_\theta^{(j)}) \delta(\sigma_\theta - \sigma_\theta^{(j)}) \delta(\sigma_y - \sigma_y^{(j)})
\end{aligned} \tag{67}$$

where $(\mu_\theta^{(j)}, \sigma_\theta^{(j)}, \sigma_y^{(j)}) \sim p(\mu_\theta, \sigma_\theta, \sigma_y | \mathcal{M}_{2b})$, $w_j \propto \prod_{i=1}^{N_D} p(D_i | \mu_\theta^{(j)}, \sigma_\theta^{(j)}, \sigma_y^{(j)}, \mathcal{M}_{2b})$, $\sum_{j=1}^{N_s} w_j = 1$

and the model evidence:

$$p(\mathcal{D} | \mathcal{M}_{2b}) \approx \frac{1}{N_s} \sum_{j=1}^{N_s} \left(\prod_{i=1}^{N_D} p(D_i | \mu_\theta^{(j)}, \sigma_\theta^{(j)}, \sigma_y^{(j)}, \mathcal{M}_{2b}) \right) \tag{68}$$

Similar to \mathcal{M}_{2a} , the marginalized posterior of θ can be estimated by:

$$p(\theta | \mathcal{D}, \mathcal{M}_{2b}) \approx \sum_{j=1}^{N_s} w_j N(\theta | \mu_\theta^{(j)}, (\sigma_\theta^{(j)})^2) \tag{69}$$

As a result, the statistics of the marginalized posterior of θ and σ_y can also be estimated based on the MCS samples:

$$\begin{aligned}
E[\theta | \mathcal{D}] &= \int \theta p(\theta | \mathcal{D}, \mathcal{M}_{2b}) d\theta \approx \sum_{j=1}^{N_s} w_j \mu_\theta^{(j)} \\
Std[\theta | \mathcal{D}] &\approx \sqrt{\left(\sum_{j=1}^{N_s} w_j \left((\sigma_\theta^{(j)})^2 + (\mu_\theta^{(j)})^2 \right) \right) - E[\theta | \mathcal{D}]^2} \\
E[\sigma_y | \mathcal{D}] &\approx \sum_{j=1}^{N_s} w_j \sigma_y^{(j)} \\
Std[\sigma_y | \mathcal{D}] &= \sqrt{\left(\sum_{j=1}^{N_s} w_j (\sigma_y^{(j)})^2 \right) - E[\sigma_y | \mathcal{D}]^2}
\end{aligned} \tag{70}$$

To perform robust prediction of a new point (\hat{x}, \hat{y}) , we need to evaluate $p((\hat{x}, \hat{y})|\mathcal{D}, \mathcal{M}_{2b})$. Again, we can use the posterior samples from Equation 67:

$$\begin{aligned} p((\hat{x}, \hat{y})|\mathcal{D}, \mathcal{M}_{2b}) &= \int p((\hat{x}, \hat{y})|\mu_\theta, \sigma_\theta, \sigma_y, \mathcal{M}_{2b}) p(\mu_\theta, \sigma_\theta, \sigma_y|\mathcal{D}, \mathcal{M}_{2b}) d\mu_\theta d\sigma_\theta d\sigma_y \\ &\approx \sum_{j=1}^{N_s} w_j p((\hat{x}, \hat{y})|\mu_\theta^{(j)}, \sigma_\theta^{(j)}, \sigma_y^{(j)}, \mathcal{M}_{2b}) \end{aligned} \quad (71)$$

Similar to \mathcal{M}_{2a} , we can evaluate $p((\hat{x}, \hat{y})|\mu_\theta^{(j)}, \sigma_\theta^{(j)}, \sigma_y^{(j)}, \mathcal{M}_{2b})$ using Equation 66 by substituting D_i with (\hat{x}, \hat{y}) .

References

- [1] K. Bae and B. K. Mallick. Gene selection using a two-level hierarchical bayesian model. *Bioinformatics*, 20(18):3423–3430, December 2004. ISSN 1367-4803. doi: 10.1093/bioinformatics/bth419.
- [2] G. Ballesteros, P. Angelikopoulos, C. Papadimitriou, and P. Koumoutsakos. Bayesian hierarchical models for uncertainty quantification in structural dynamics. In M. Beer, S.K. Au, and J. W. Hall, editors, *Vulnerability, Uncertainty, and Risk: Quantification, Mitigation, and Management*, volume 162, pages 1615–1624. American Society of Civil Engineers (ASCE), Reston, Virginia, USA, 2014.
- [3] M. Barrault, Y. Maday, N. C. Nguyen, and A. T. Patera. An ‘empirical interpolation’ method: application to efficient reduced-basis discretization of partial differential equations. *C. R. Math.*, 339(9):667–672, November 2004.
- [4] J.L. Beck. Bayesian system identification based on probability logic. *Structural Control and Health Monitoring*, 17(7):825–847, 2010.
- [5] J.L. Beck and K.V. Yuen. Model selection using response measurements: Bayesian probabilistic approach. *J. Eng. Mech.-ASCE*, 130(2):192–203, February 2004. ISSN 0733-9399. doi: 10.1061/(ASCE)0733-9399(2004)130:2(192).
- [6] D. Calvetti and E. Somersalo. Hypermodels in the bayesian imaging framework. *Inverse Probl.*, 24(3), June 2008. ISSN 0266-5611. doi: 10.1088/0266-5611/24/3/034013.
- [7] J.Y. Ching and Y.C. Chen. Transitional markov chain monte carlo method for bayesian model updating, model class selection, and model averaging. *Journal of Engineering Mechanics-ASCE*, 133(7):816–832, July 2007.
- [8] P.D. Congdon. *Applied Bayesian Hierarchical Methods*. CRC Press, 2010.

- [9] L. Fei-Fei and P. Perona. A bayesian hierarchical model for learning natural scene categories. In *Proceedings – IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 524–531, San Diego, CA, USA, 2005.
- [10] S.D. Finley, P. Angelikopoulos, P. Koumoutsakos, and A.S. Popel. Pharmacokinetics of Anti-VEGF agent aflibercept in cancer predicted by data-driven, molecular-detailed model. *CPT: Pharmacometrics & Systems Pharmacology*, 4:641–649, 2015.
- [11] A. Gelman and J. Hill. *Data Analysis Using Regression and Multi-level/Hierarchical Models*. Cambridge University Press, 2006.
- [12] A. Gelman, J.B. Carlin, H.S. Stern, and D.B. Rubin. *Bayesian Data Analysis*. Chapman & Hall/CRC, 2 edition, 2004.
- [13] N. Guha, X. Wu, Y. Efendiev, B. Jin, and B. K. Mallick. A variational bayesian approach for inverse problems with skew-t error distributions. *J. Comput. Phys.*, 301:377–393, November 2015. ISSN 0021-9991. doi: 10.1016/j.jcp.2015.07.062.
- [14] P. Richard Hahn, Indranil Goswami, and Carl F. Mela. A bayesian hierarchical model for inferring player strategy types in a number guessing game. *Ann. Appl. Stat.*, 9(3):1459–1483, September 2015. ISSN 1932-6157. doi: 10.1214/15-AOAS830.
- [15] J. S. Hesthaven, B. Stamm, and S. Zhang. Efficient greedy algorithms for high-dimensional parameter spaces with applications to empirical interpolation and reduced basis methods. *ESAIM-Math. Model. Numer. Anal.-Model. Math. Anal. Numer.*, 48(1):259–283, January 2014. ISSN 0764-583X. doi: 10.1051/m2an/2013100.
- [16] B.M. Hill. Inference about variance components in the one-way model. *J. Am. Stat. Assoc.*, 60(311):806–825, 1965. ISSN 0162-1459.
- [17] Y. Huang, J.L. Beck, S. Wu, and H. Li. Robust bayesian compressive sensing for signals in structural health monitoring. *Computer-Aided Civil and Infrastructure Engineering*, 29(3):160–179, 2014.
- [18] S. Ji, Y. Xue, and L. Carin. Bayesian compressive sensing. *IEEE Trans. Signal Process.*, 56(6):2346–2356, June 2008. ISSN 1053-587X. doi: 10.1109/TSP.2007.914345.
- [19] D. Koller and N. Friedman. *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, 2009.
- [20] D.J.C. Mackay. Bayesian non-linear modelling for the prediction competition. *ASHRAE Trans.*, 100(2):1053–1062, 1994.

- [21] J. B. Nagel and B. Sudret. A unified framework for multilevel uncertainty quantification in bayesian inverse problems. *Probabilistic Engineering Mechanics*, [Accepted], 2015.
- [22] K. Sargsyan, H. N. Najm, and R. Ghanem. On the statistical calibration of physical models. *Int. J. Chem. Kinet.*, 47(4):246–276, April 2015. ISSN 0538-8066. doi: 10.1002/kin.20906.
- [23] M. Sato, T. Yoshioka, S. Kajihara, K. Toyama, N. Goda, K. Doya, and M Kawato. Hierarchical bayesian estimation for meg inverse problem. *Neuroimage*, 23(3):806–826, 2004.
- [24] G.C. Tiao and W.Y. Tan. Bayesian analysis of random-effect models in analysis of variance. I. posterior distribution of variance-components. *Biometrika*, 52(1/2):37–53, 1965. ISSN 0006-3444.
- [25] M.E. Tipping. Sparse bayesian learning and the relevance vector machine. *Journal of Machine Learning Research*, 1(3):211–244, 2001.
- [26] M.E. Tipping. Bayesian inference: An introduction to principles and practice in machine learning. In *Advanced Lectures on Machine Learning*, volume 3176, pages 41–62. Springer, 2004.
- [27] M.E. Tipping and A.C. Faul. Fast marginal likelihood maximization for sparse bayesian models. In C.M. Bishop and B.J. Frey, editors, *Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics*, Key West, FL, January 2003.
- [28] M. W. Woolrich, T. E. J. Behrens, C. F. Beckmann, M. Jenkinson, and S. M. Smith. Multilevel linear modelling for FMRI group analysis using bayesian inference. *Neuroimage*, 21(4):1732–1747, April 2004. ISSN 1053-8119. doi: 10.1016/j.neuroimage.2003.12.023.
- [29] S. Wu, P. Angelikopoulos, C. Papadimitriou, R. Moser, and P. Koumoutsakos. A hierarchical bayesian framework for force field selection in molecular dynamics simulations. *Phil. Trans. R. Soc. A*, 374(2060):20150032, 2015.
- [30] S. Wu, P. Angelikopoulos, G. Tauriello, and P. Koumoutsakos. Hierarchical bayesian modeling for molecular dynamics simulations using heterogeneous data. In *Proceedings of UNCECOMP*, 2015.
- [31] S. Wu, P. Angelikopoulos, C. Papadimitriou, and P. Koumoutsakos. Bayesian annealed sequential importance sampling (BASIS): an unbiased version of transitional markov chain monte carlo. *ASCE-ASME Journal of Risk and Uncertainty in Engineering Systems, Part B: Mechanical Engineering*, [Under Review], 2016.